

P&S Heterogeneous Systems

Collaborative Computing

Dr. Juan Gómez Luna

Prof. Onur Mutlu

ETH Zürich

Fall 2022

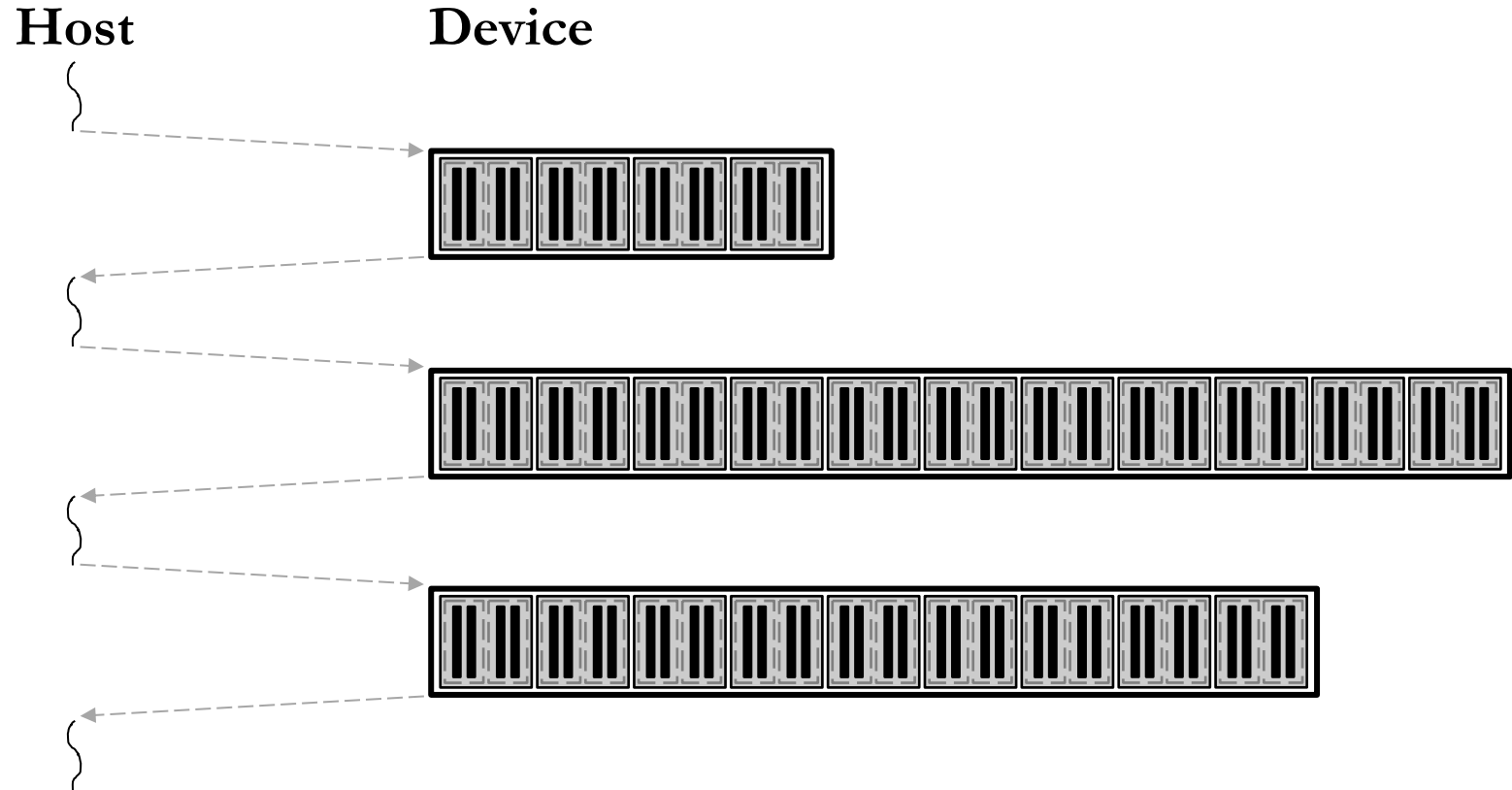
16 January 2023

In Our Previous Lecture...

Dynamic Parallelism

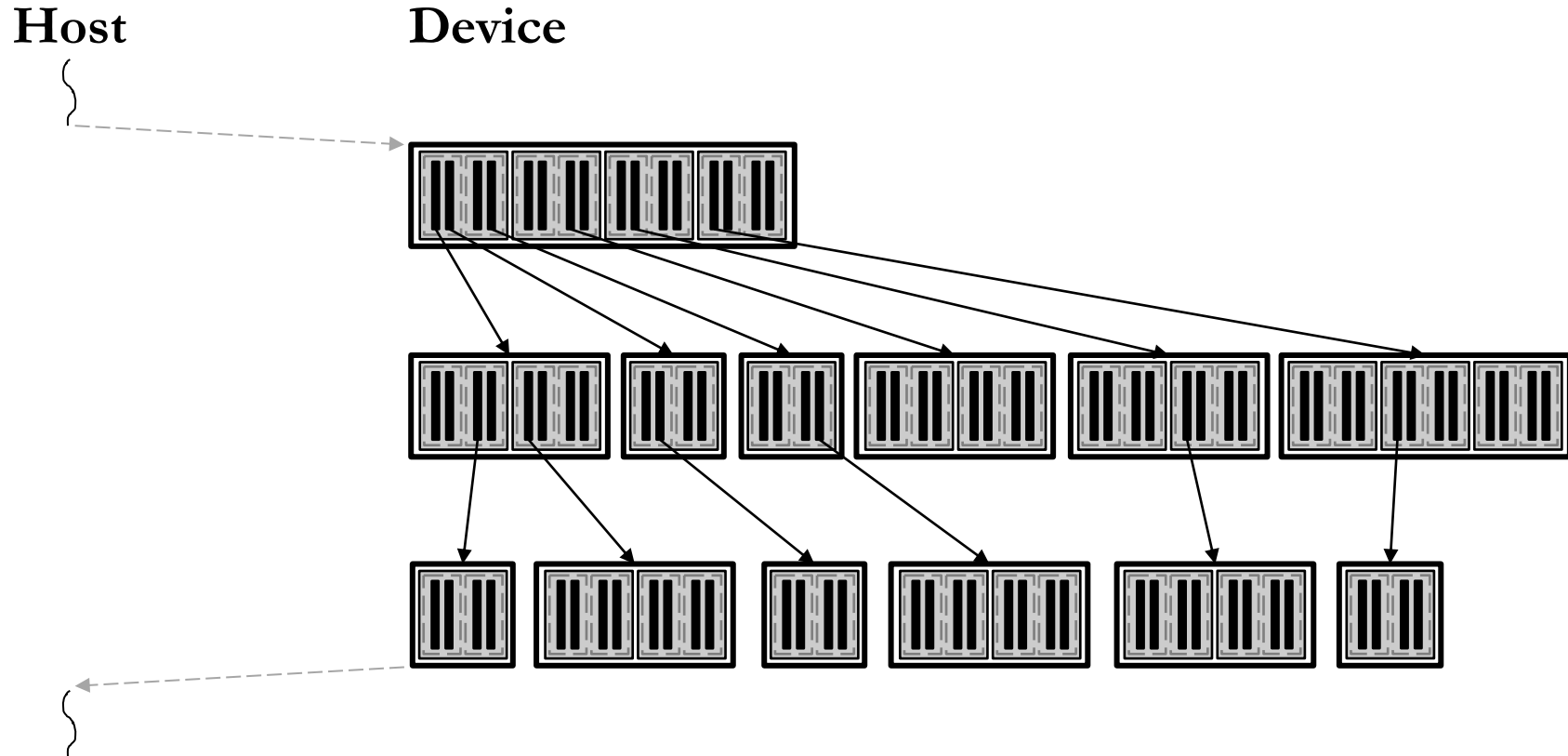
- GPU programming frameworks provide an interface to express **dynamic refinement algorithms** in a more natural way
 - This dynamic parallelism interface allows GPU threads to launch GPU kernels when new work is dynamically discovered
 - Recall BFS
 - Each node in the frontier has a different number of neighbors
- CUDA Dynamic Parallelism
 - Important semantics when a kernel is launched from a kernel
 - Performance considerations

Kernel Launch without Dynamic Parallelism



Previously, kernels could **only be launched from the host** (painful to program!)

Kernel Launch with Dynamic Parallelism

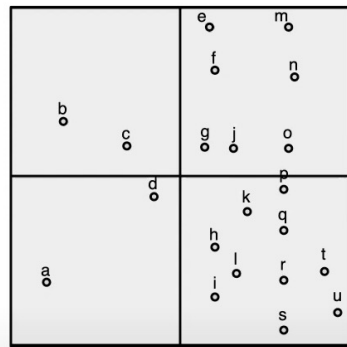


Easier to write programs with **dynamically discovered parallelism**

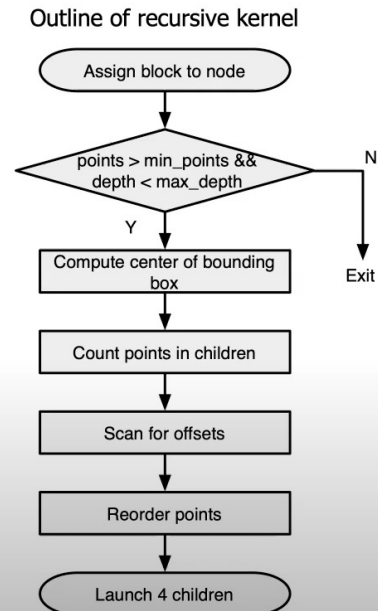
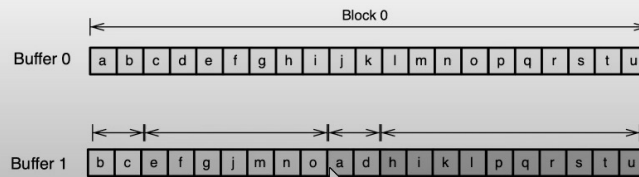
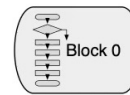
Lecture on Dynamic Parallelism

A Recursive Example: Quadtree (V)

- 1 thread block is launched from host



Depth = 0



HetSys Course: Lecture 13: Dynamic Parallelism (Spring 2022)

247 views • Premiered Jun 7, 2022

12 DISLIKE SHARE CLIP SAVE ...



Onur Mutlu Lectures
25.6K subscribers

SUBSCRIBED



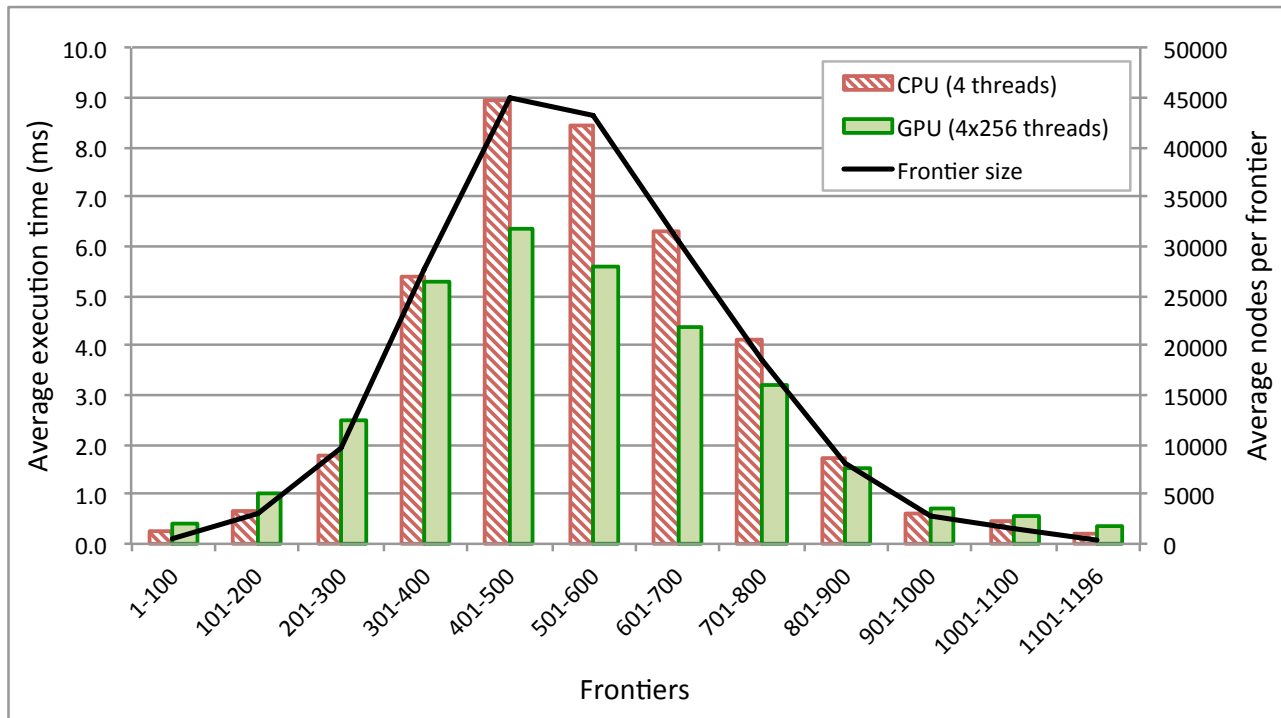
<https://youtu.be/X74BLPO8tT4>

Collaborative Computing

Recall: BFS on CPU or GPU?

■ Motivation

- Small-sized frontiers underutilize GPU resources
 - NVIDIA Jetson TX1 (4 ARMv8 CPU cores + 2 GPU cores)
 - New York City roads



BFS: Collaborative Implementation (I)

- Choose CPU or GPU depending on frontier

```
// Host code
while(frontier_size != 0){

    if(frontier_size < LIMIT){

        // Launch CPU threads
    }
    else{

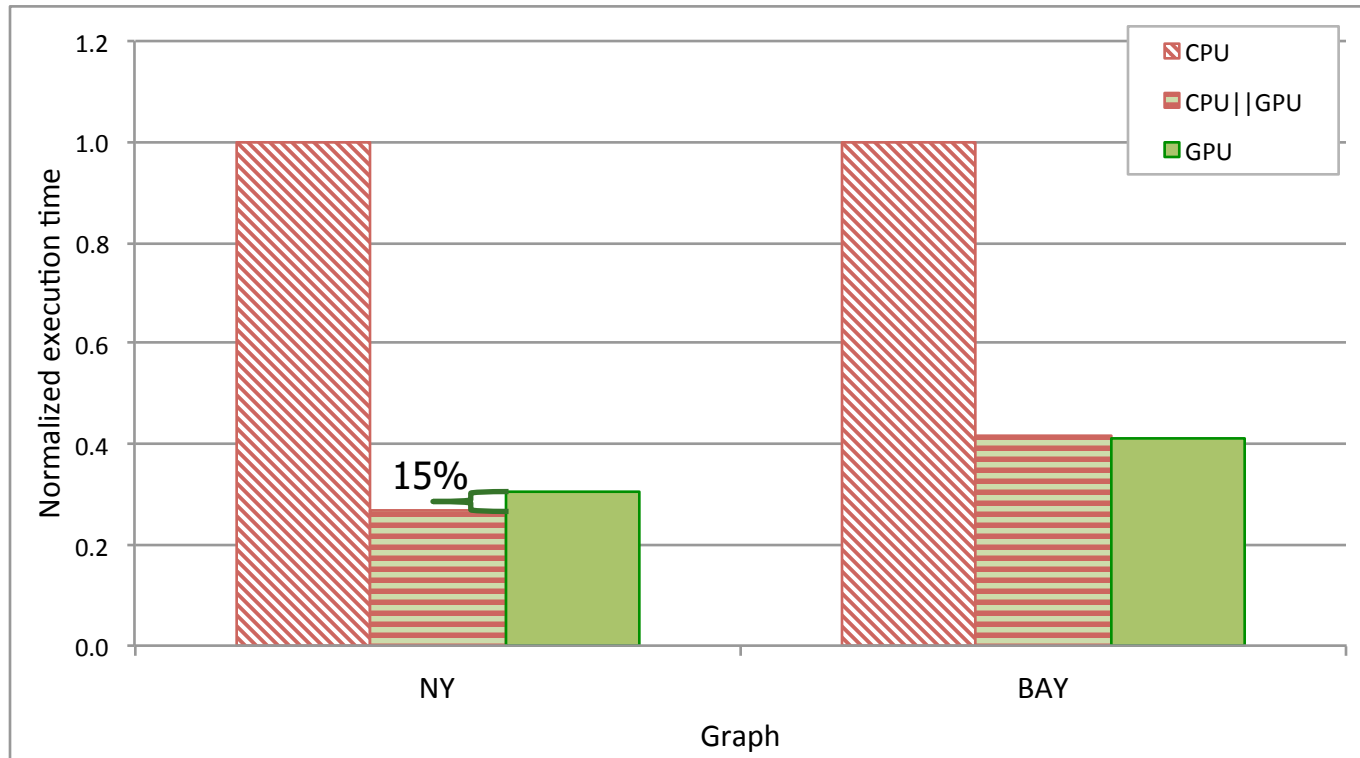
        // Launch GPU kernel
    }
}
```

- CPU threads or GPU kernel keep running while the condition is satisfied

BFS: Collaborative Implementation (II)

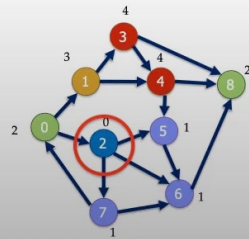
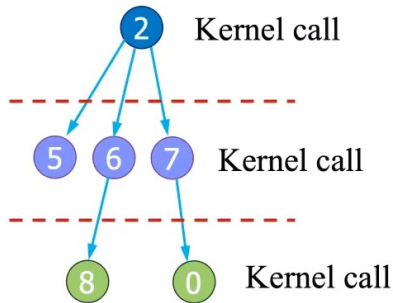
■ Experimental results

- NVIDIA Jetson TX1 (4 ARMv8 CPU cores + 2 GPU cores)



Lecture on Graph Search

Kernel Arrangement



- Creating global barriers needs frequent kernel launches
- Too much overhead
- Solutions:
 - ❑ Partially use GPU-synchronization
 - ❑ Multi-layer Kernel Arrangement
 - ❑ Dynamic Parallelism
 - ❑ Persistent threads with global barriers



HetSys Course: Lecture 11: Parallel Patterns: Graph Search (Spring 2022)

325 views • Premiered May 24, 2022

15 DISLIKE SHARE CLIP SAVE ...



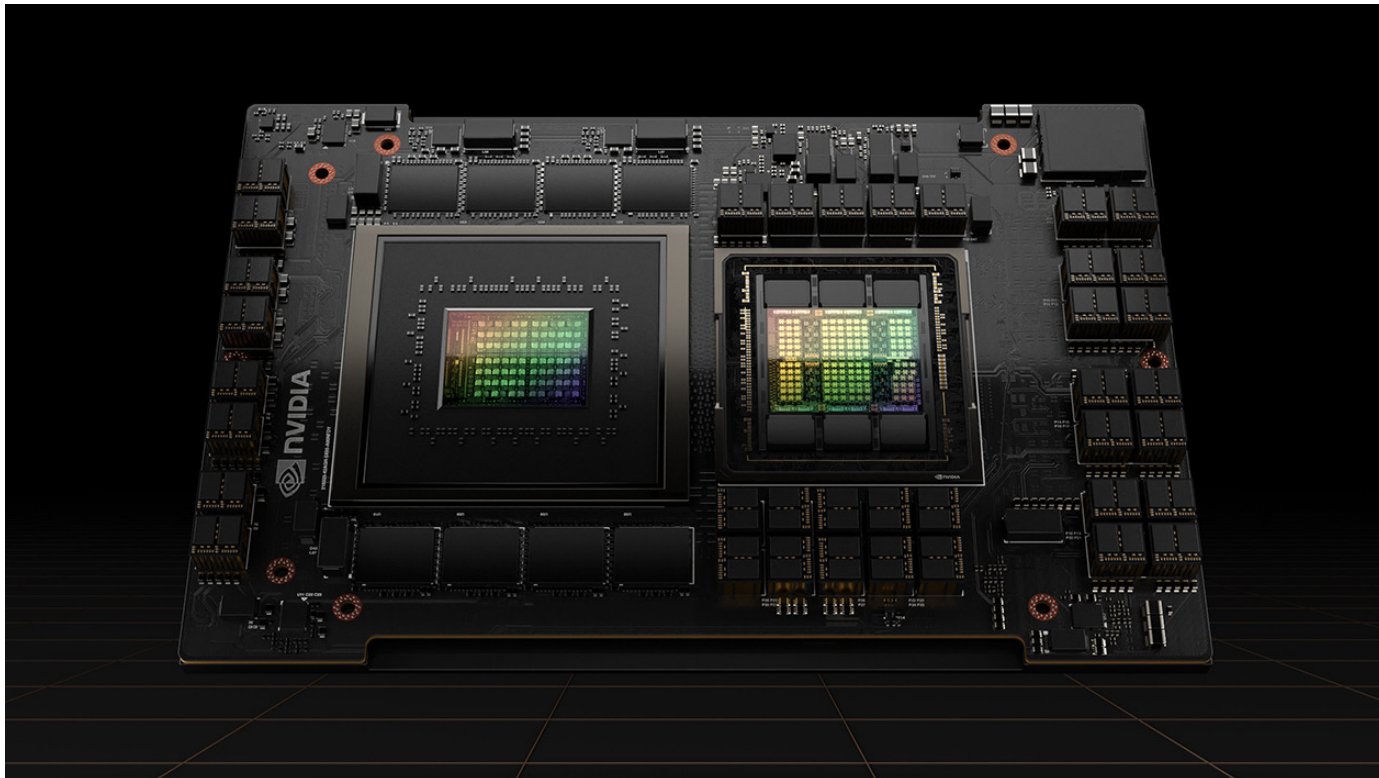
Onur Mutlu Lectures
25.6K subscribers

SUBSCRIBED



NVIDIA Grace Hopper Superchip

- CPU + GPU
 - Grace CPU + Hopper GPU
- 900 GB/s coherent interface (7x faster than PCIe Gen 5)



Unified Memory

Memory Allocation and Data Transfers

- Traditional approach to **device allocation, CPU-GPU transfer, and GPU-CPU transfer**
 - ❑ `cudaMalloc()`;
 - ❑ `cudaMemcpy()`;
- Naturally matches systems with **discrete GPUs**

```
// Allocate input
malloc(input, ...);
cudaMalloc(d_input, ...);
cudaMemcpy(d_input, input, ..., HostToDevice); // Copy to device memory

// Allocate output
malloc(output, ...);
cudaMalloc(d_output, ...);

// Launch GPU kernel
gpu_kernel<<<blocks, threads>>> (d_output, d_input, ...);

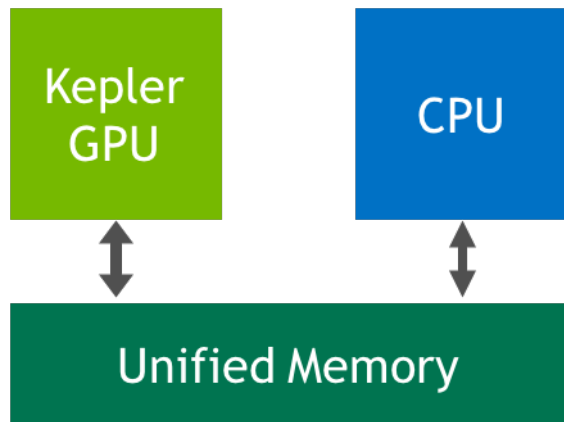
// Synchronize
cudaDeviceSynchronize();

// Copy output to host memory
cudaMemcpy(output, d_output, ..., DeviceToHost);
```

Unified Memory

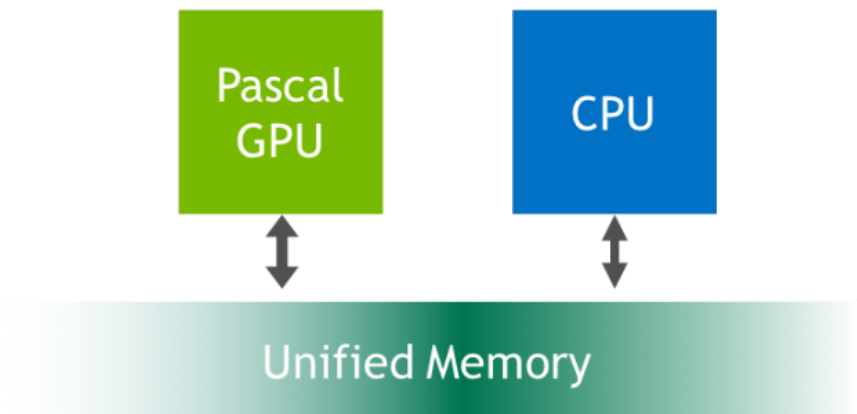
- Unified Virtual Address space
 - Same virtual address space across host and device
- CUDA 6.0: Unified memory
- CUDA 8.0 + Pascal: GPU page faults

CUDA 6 Unified Memory



(Limited to GPU Memory Size)

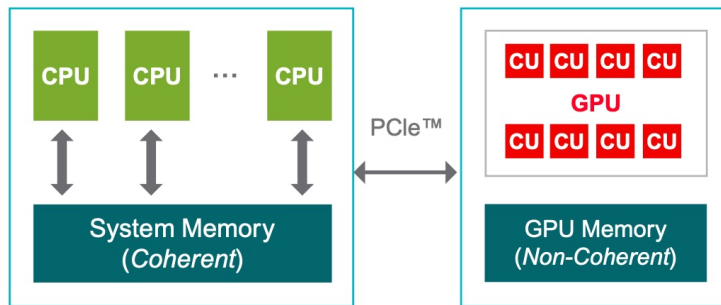
Pascal Unified Memory



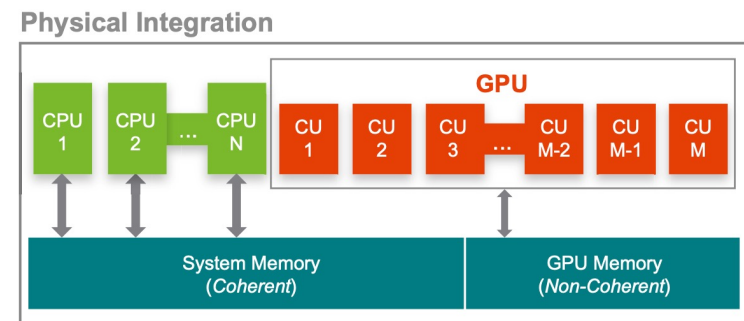
(Limited to System Memory Size)

Heterogeneous System Architecture

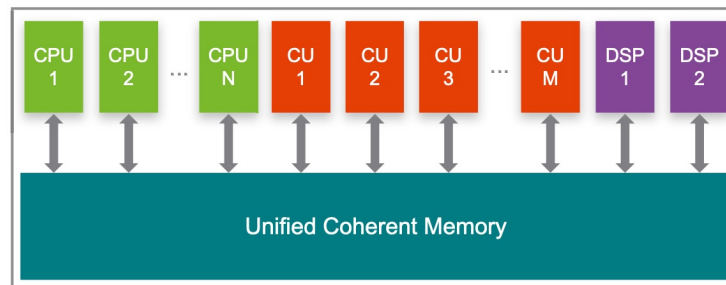
- HSA extends the unified memory space beyond GPUs
 - DSPs, DMA engines, cryptoengines, and other accelerators



Legacy GPU compute on discrete GPU cards



Legacy GPU compute on SOC's



An HSA enabled SOC featuring multiple processors beyond CPU

Unified Memory: Memory Management

- Easier programming with Unified Memory

- ❑ `cudaMallocManaged()`;

```
// Allocate input
malloc(input, ...);
cudaMallocManaged(d_input, ...);
memcpy(d_input, input, ...); // Copy to managed memory

// Allocate output
cudaMallocManaged(d_output, ...);

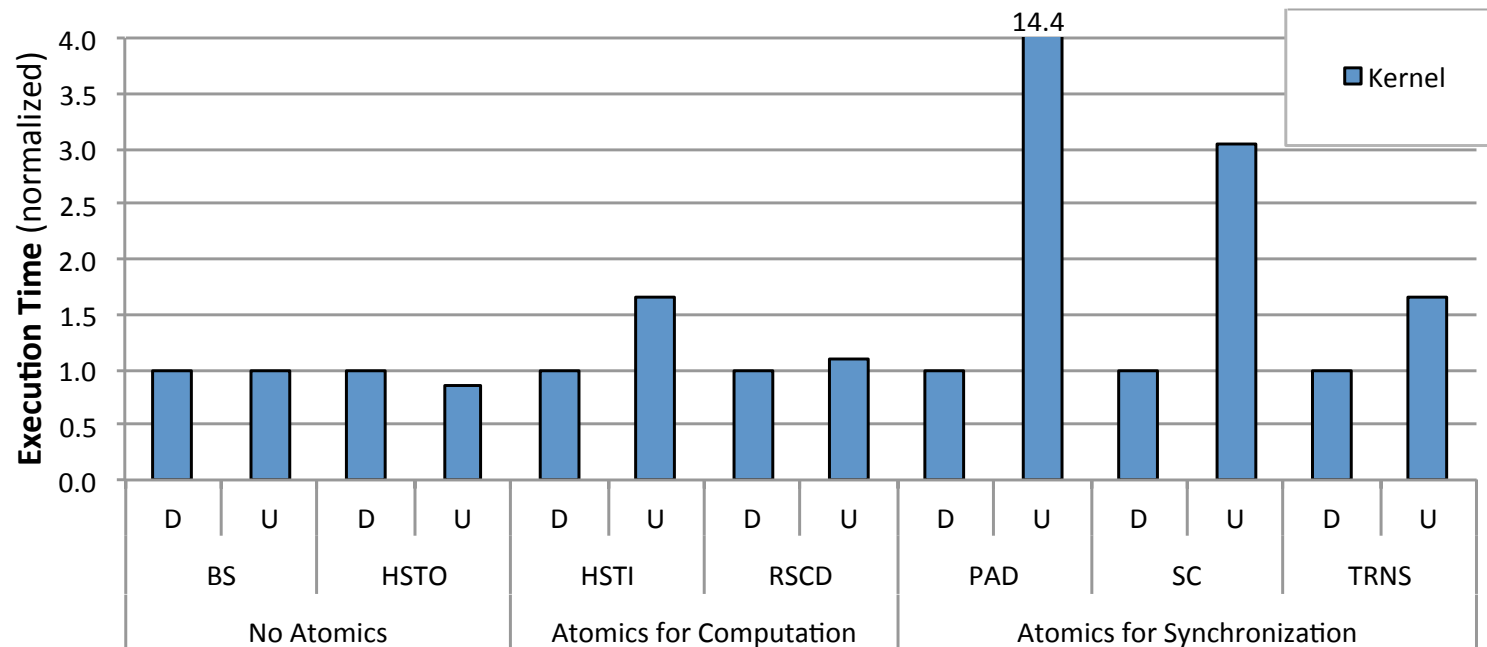
// Launch GPU kernel
gpu_kernel<<<blocks, threads>>> (d_output, d_input, ...);

// Synchronize
cudaDeviceSynchronize();
```

- No need for double allocation or explicit data transfers
- Naturally matches physically integrated devices (e.g., CPU and GPU in the same chip) or devices with the same physical memory (e.g., CPU and GPU in the same package)
 - ❑ But it can also be implemented for discrete GPUs

Unified Memory: Kernel Time

- IBM Power8 with NVIDIA Pascal GPU
 - **D**: Discrete (or traditional, without unified memory)
 - **U**: Unified memory

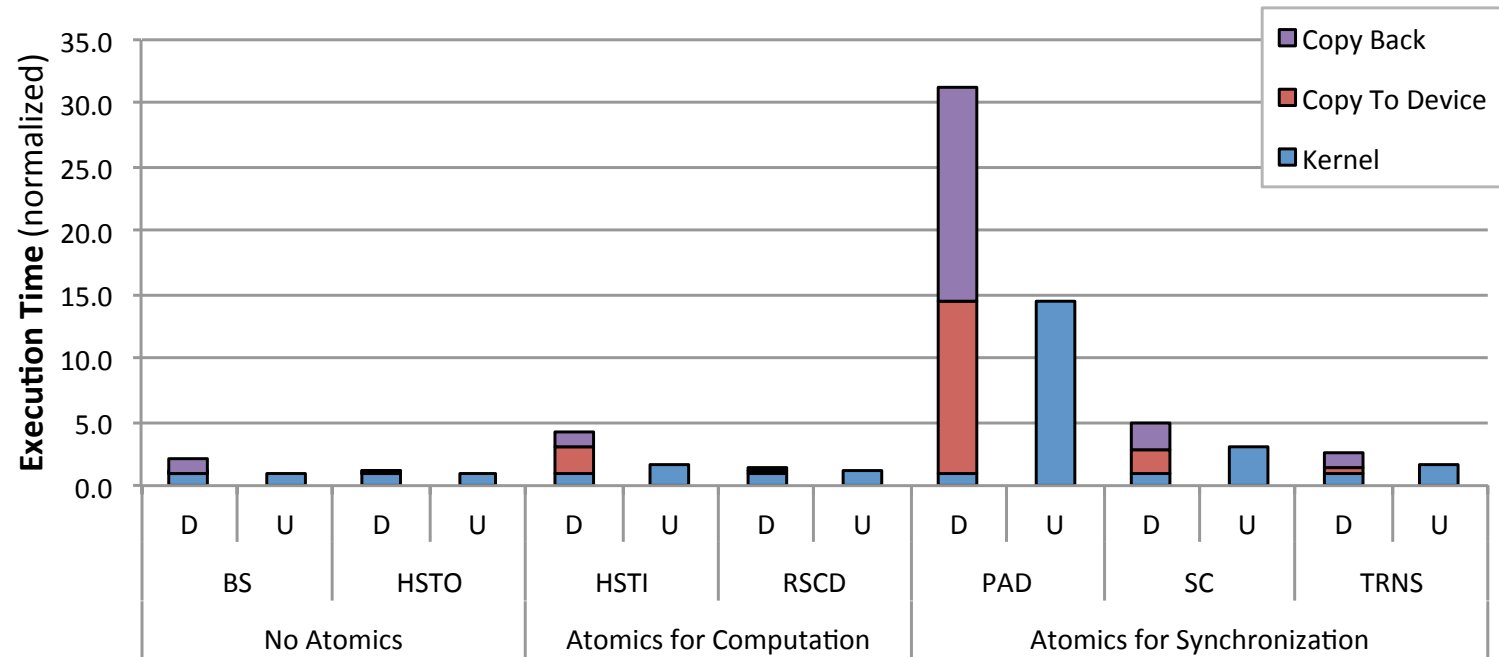


No cross-device
communication

Cross-device communication may heavily
burden kernel performance

Unified Memory: Total Execution Time

- IBM Power8 with NVIDIA Pascal GPU
 - **D**: Discrete (or traditional, without unified memory)
 - **U**: Unified memory



Unified memory can hide data transfers with kernel execution

How to Implement Collaborative Computing Applications?

Collaborative Computing Applications

- Case studies using CPU and GPU
- Kernel launches are asynchronous
 - CPU can work while waits for GPU to finish
 - Traditionally, this is the most efficient way to exploit heterogeneity

```
// Allocate input
malloc(input, ...);
cudaMalloc(d_input, ...);
cudaMemcpy(d_input, input, ..., HostToDevice); // Copy to device memory

// Allocate output
malloc(output, ...);
cudaMalloc(d_output, ...);

// Launch GPU kernel
gpu_kernel<<<blocks, threads>>> (d_output, d_input, ...);

// CPU can do things here

// Synchronize
cudaDeviceSynchronize();

// Copy output to host memory
cudaMemcpy(output, d_output, ..., DeviceToHost);
```

Fine-Grained Collaboration

- Fine-grained collaboration becomes possible with unified memory (post Kepler/Maxwell architecture)
- Pascal/Volta/Turing/Ampere Unified Memory (& HSA)
 - CPU-GPU memory coherence
 - System-wide atomic operations

```
// Allocate input
cudaMallocManaged(input, ...);

// Allocate output
cudaMallocManaged(output, ...);

// Launch GPU kernel
gpu_kernel<<<blocks, threads>>> (output, input, ...);

// CPU can do things here
output[x] = input[y];

output[x+1].fetch_add(1);
```

CUDA 8.0 and Later

- Unified memory

```
cudaMallocManaged(&h_in, in_size);
```

- System-wide atomics

```
old = atomicAdd_system(&h_out[x], inc);
```

OpenCL 2.0 and Later

■ Shared virtual memory

```
XYZ * h_in = (XYZ *)clSVMAlloc(  
    ocl.clContext, CL_MEM_SVM_FINE_GRAIN_BUFFER, in_size, 0);
```

■ More flags:

```
CL_MEM_READ_WRITE  
CL_MEM_SVM_ATOMICS
```

■ C++11 atomic operations

(memory_scope_all_svm_devices)

```
old = atomic_fetch_add(&h_out[x], inc);
```

C++AMP (HCC)

- Unified memory space (HSA)

```
XYZ *h_in = (XYZ *)malloc(in_size);
```

- C++11 atomic operations

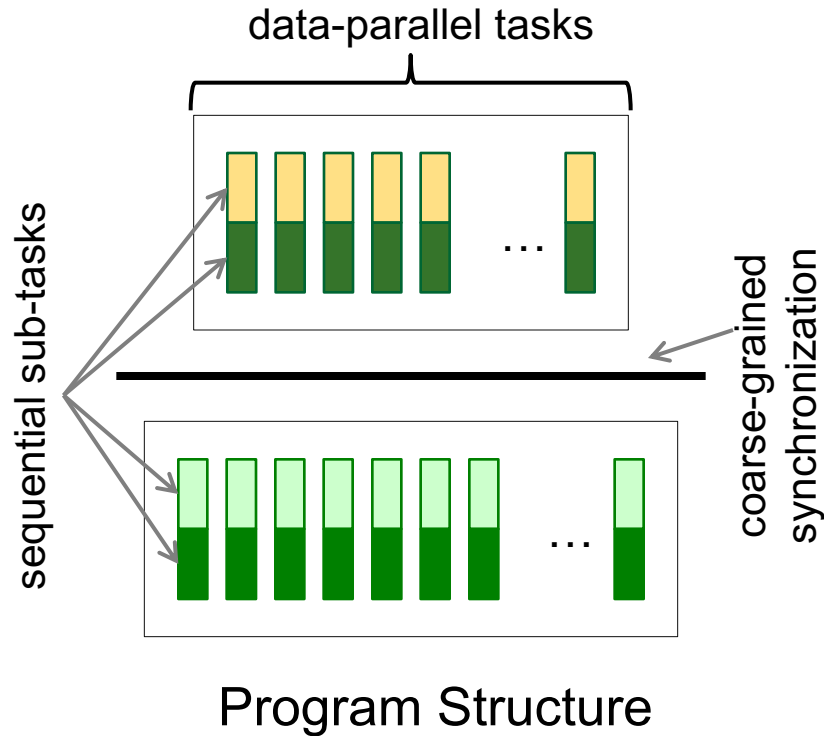
```
(memory_scope_all_svm_devices)
```

- Platform atomics (HSA)

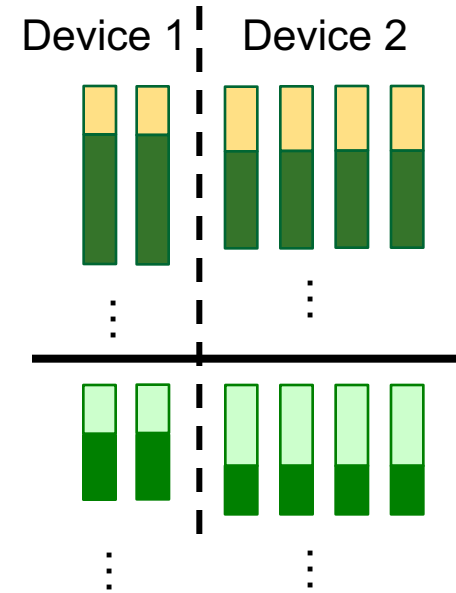
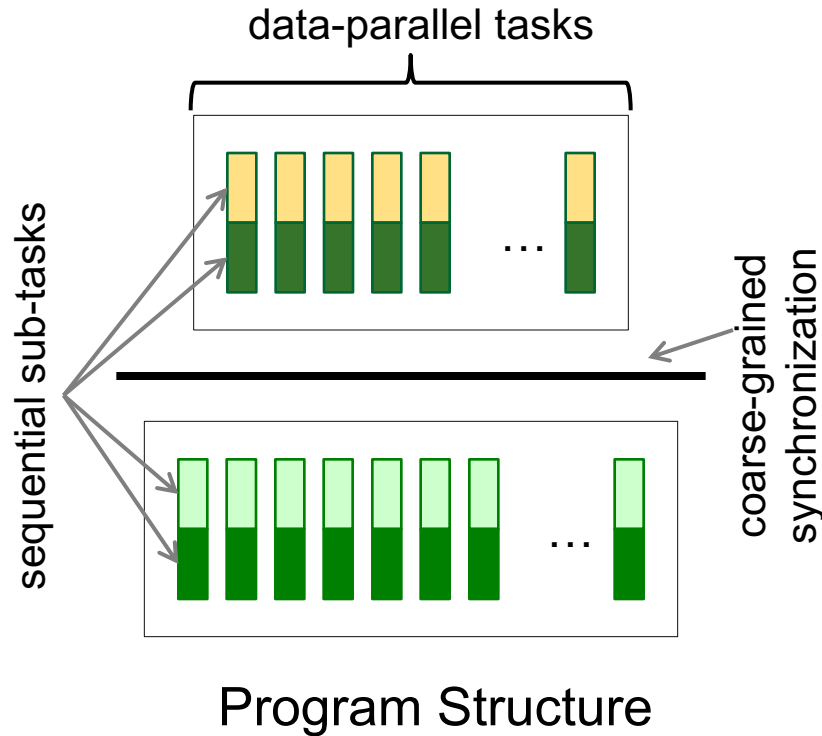
```
old = atomic_fetch_add(&h_out[x], inc);
```

Collaborative Patterns

Traditional Program Structure

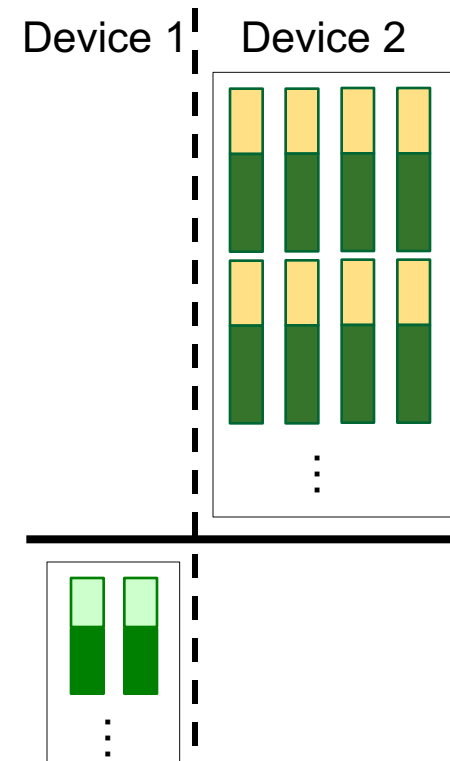
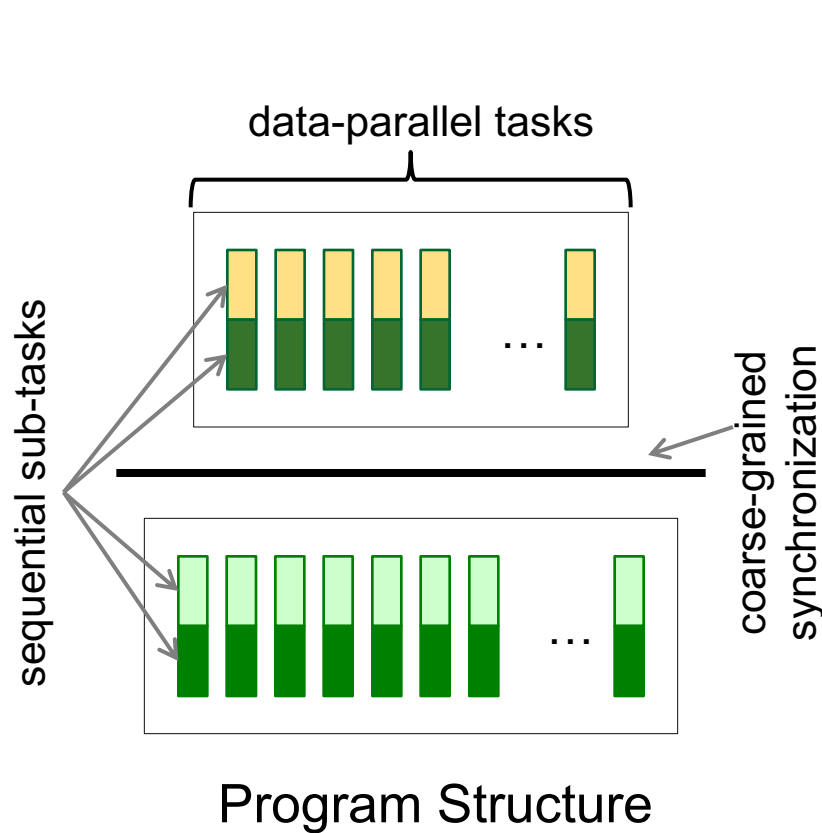


Collaborative Patterns: Data Partitioning



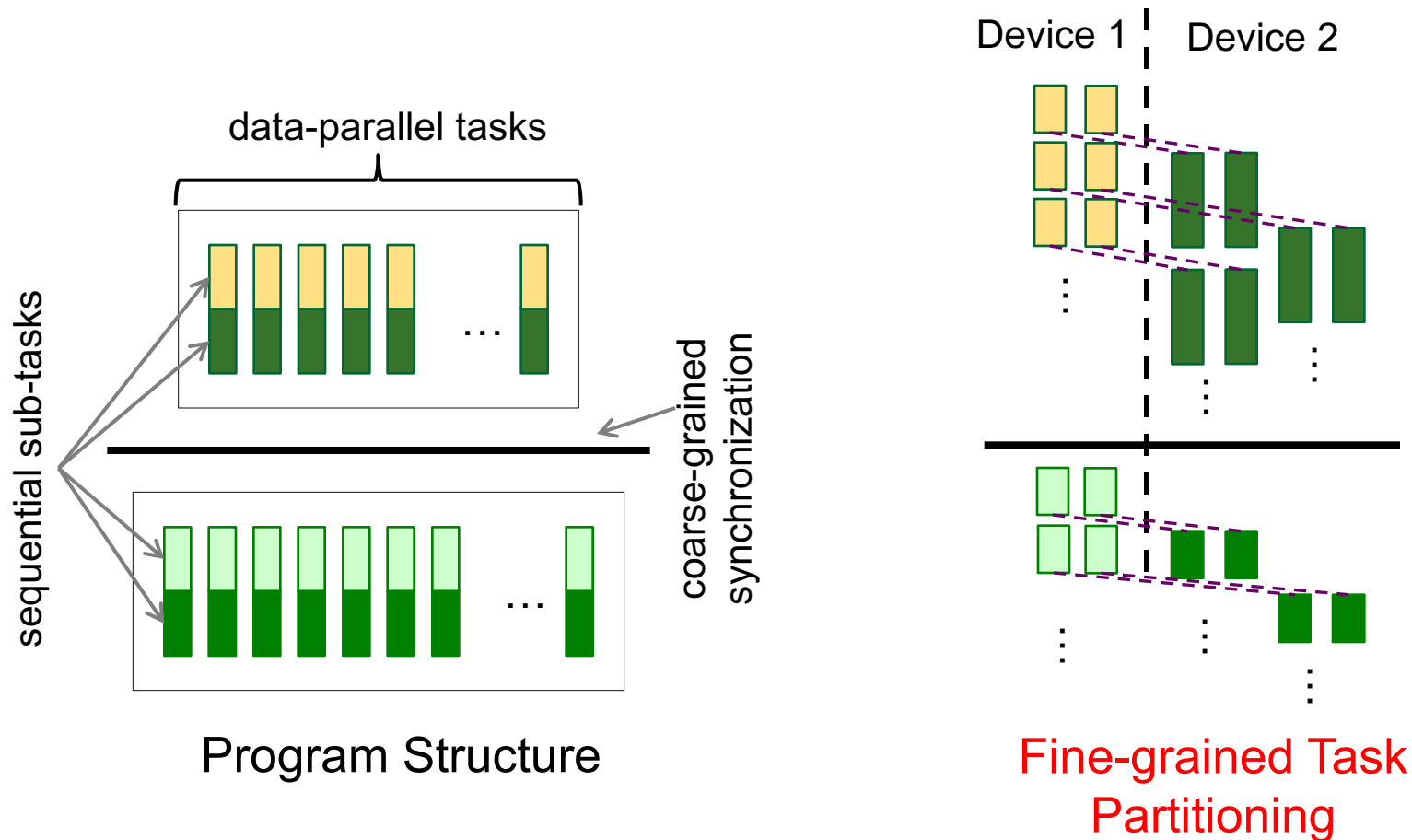
Data Partitioning

Collaborative Patterns: Task Partitioning (I)



Coarse-grained Task Partitioning

Collaborative Patterns: Task Partitioning (II)

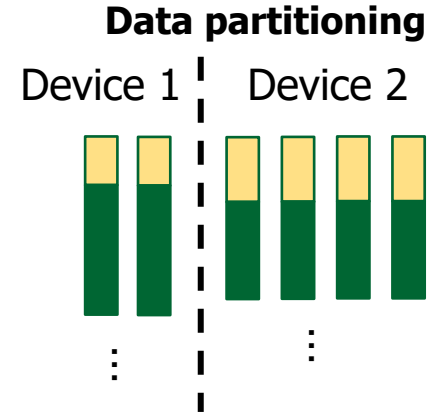


Analytical Modeling

- N : Number of data parallel tasks in the application
- $t_{i,D1}$: Execution time of sub-task i by a Device 1 worker
- $t_{i,D2}$: Execution time of sub-task i by a Device 2 worker
- w_{D1} : Number of available Device 1 workers
- w_{D2} : Number of available Device 2 workers
- β : Distribution and aggregation overhead factor
- α : Fraction of data parallel tasks assigned to Device 1
- S_{D1} and S_{D2} are, respectively, the set of subtasks/tasks executed in Device 1 and Device 2

Analytical Model: Data Partitioning

- N : Number of data parallel tasks in the application
- $t_{i,D1}$: Execution time of sub-task i by a Device 1 worker
- $t_{i,D2}$: Execution time of sub-task i by a Device 2 worker
- w_{D1} : Number of available Device 1 workers
- w_{D2} : Number of available Device 2 workers
- β : Distribution and aggregation overhead factor
- α : Fraction of data parallel tasks assigned to Device 1



Data partitioning

The total execution time is

$$t_{\text{data, total}} = \beta_{\text{data}} \cdot \max \left(\frac{\alpha N \sum_i t_{i,D1}}{w_{D1}}, \frac{(1 - \alpha) N \sum_i t_{i,D2}}{w_{D2}} \right)$$

Total D1 execution time (sequential execution) Total D2 execution time (sequential execution)

Fixing all the variables except α , the optimal α (global minimum point) is

$$\alpha^* = \frac{\sum_i t_{i,D2}}{w_{D2}} / \left(\frac{\sum_i t_{i,D1}}{w_{D1}} + \frac{\sum_i t_{i,D2}}{w_{D2}} \right)$$

Workloads of Device 1 and Device 2 workers are balanced

Analytical Model: Fine-Grained Task Part.

- N : Number of data parallel tasks in the application
- $t_{i,D1}$: Execution time of sub-task i by a Device 1 worker
- $t_{i,D2}$: Execution time of sub-task i by a Device 2 worker
- w_{D1} : Number of available Device 1 workers
- w_{D2} : Number of available Device 2 workers
- β : Distribution and aggregation overhead factor
- S_{D1} and S_{D2} are, respectively, the set of subtasks executed in Device 1 and Device 2

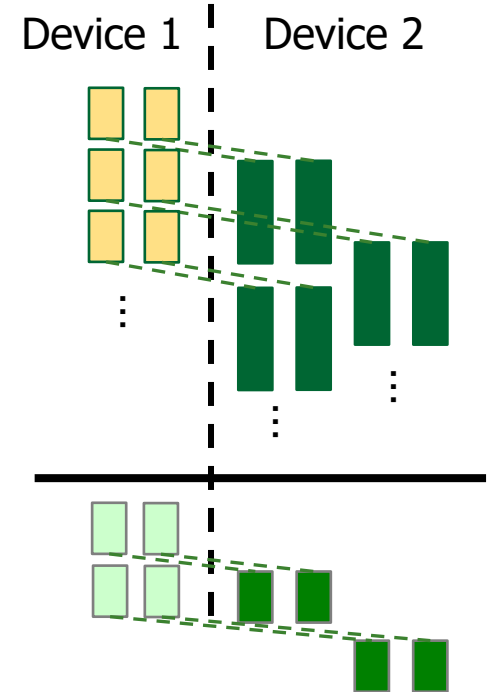
Fine-grained task partitioning

The total execution time is

$$t_{\text{task, total}} = \beta_{\text{task}} N \cdot \max \left(\frac{\sum_{i \in S_{D1}} t_{i,D1}}{w_{D1}}, \frac{\sum_{i \in S_{D2}} t_{i,D2}}{w_{D2}} \right)$$

(Assume sub-tasks are very fine-grained)

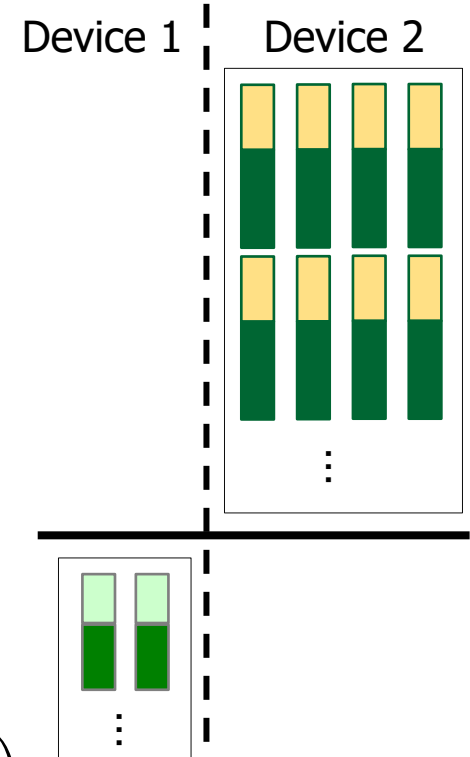
Fine-grained task partitioning



Analytical Model: Coarse-Grained Task Part.

- N : Number of data parallel tasks in the application
- $t_{i,D1}$: Execution time of sub-task i by a Device 1 worker
- $t_{i,D2}$: Execution time of sub-task i by a Device 2 worker
- w_{D1} : Number of available Device 1 workers
- w_{D2} : Number of available Device 2 workers
- β : Distribution and aggregation overhead factor
- S_{D1} and S_{D2} are, respectively, the set of tasks executed in Device 1 and Device 2

Coarse-grained task partitioning



Coarse-grained task partitioning

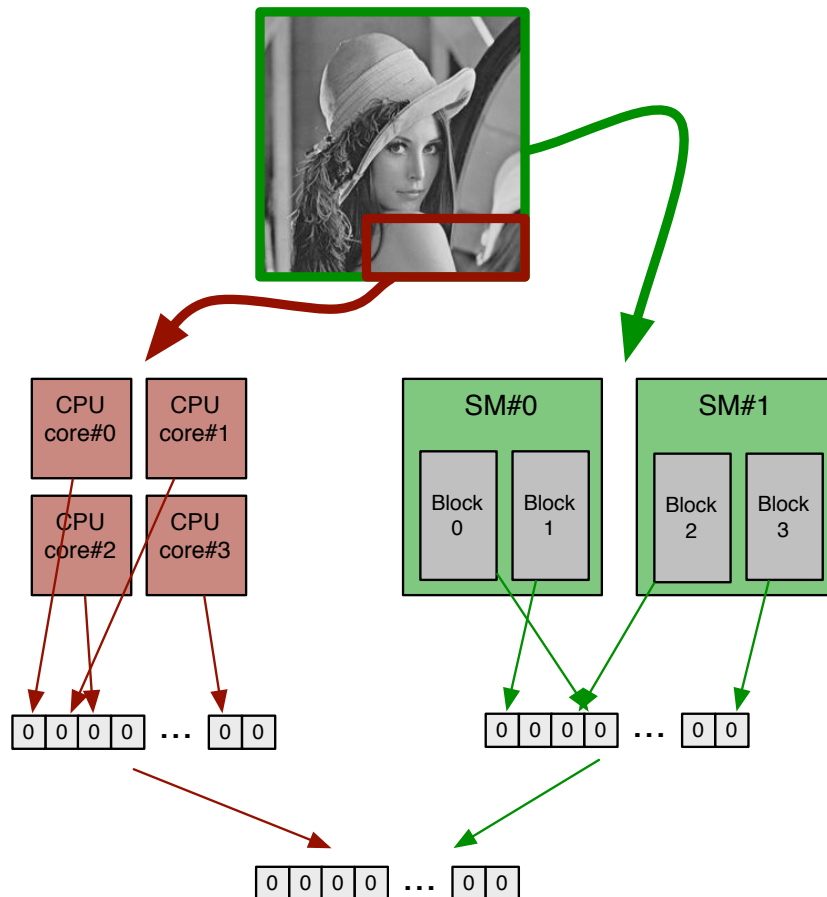
The total execution time is

$$t_{\text{task, total}} = \beta_{\text{task}} N \cdot \left(\frac{\sum_{i \in S_{D1}} t_{i,D1}}{w_{D1}} + \frac{\sum_{i \in S_{D2}} t_{i,D2}}{w_{D2}} \right)$$

Data Partitioning

Histogram without Unified Memory

- Traditional approach: **Separate CPU and GPU histograms** are merged at the end



```
malloc(CPU image);
cudaMalloc(GPU image);
cudaMemcpy(GPU image, CPU image, ...,
           HostToDevice);
malloc(CPU histogram);
memset(CPU histogram, 0);
cudaMalloc(GPU histogram);
cudaMemset(GPU histogram, 0);

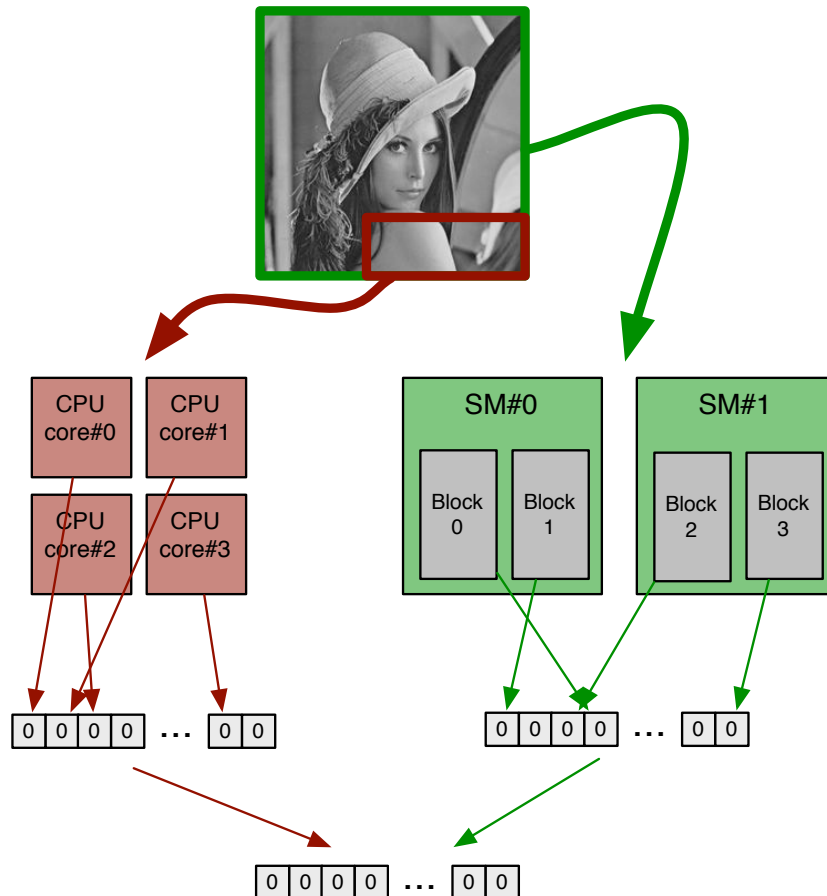
// Launch CPU threads
// Launch GPU kernel

cudaMemcpy(GPU histogram, DeviceToHost);

// Launch CPU threads for merging
```


Histogram with Unified Memory (I)

- Traditional approach: **Separate CPU and GPU histograms** are merged at the end



```
malloc(CPU image);  
cudaMallocManaged(GPU image);  
memcpy(GPU image, CPU image, ...);
```

```
malloc(CPU histogram);  
memset(CPU histogram, 0);  
cudaMallocManaged(GPU histogram);  
cudaMemset(GPU histogram, 0);
```

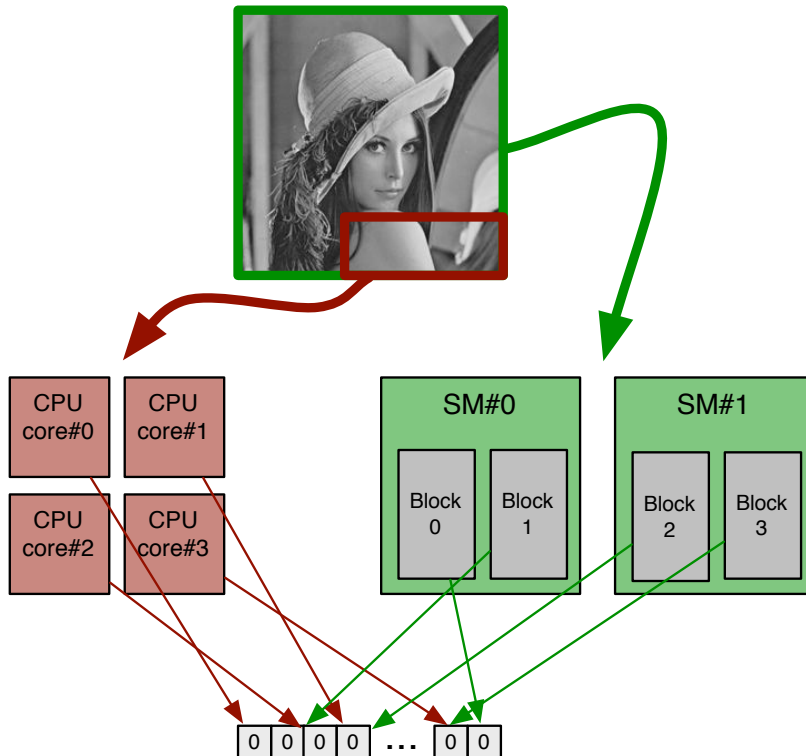
```
// Launch CPU threads  
// Launch GPU kernel
```

```
cudaDeviceSynchronize();
```

```
// Launch CPU threads for merging
```

Histogram with Unified Memory (II)

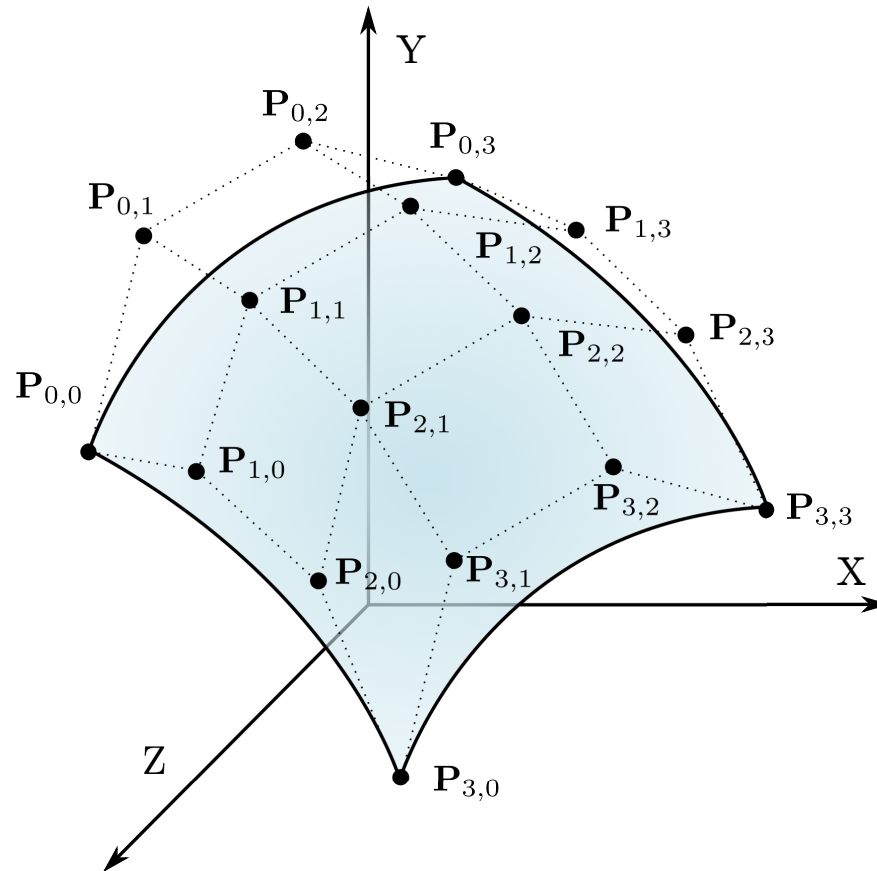
- System-wide atomic operations: **One single histogram**



```
malloc(CPU image);  
cudaMallocManaged(GPU image);  
memcpy(GPU image, CPU image, ...);  
  
cudaMallocManaged(Histogram);  
cudaMemset(Histogram, 0);  
  
// Launch CPU threads  
// Launch GPU kernel (atomicAdd_system)
```

Bézier Surfaces (I)

- Bézier surface: 4x4 net of control points



Bézier Surfaces (II)

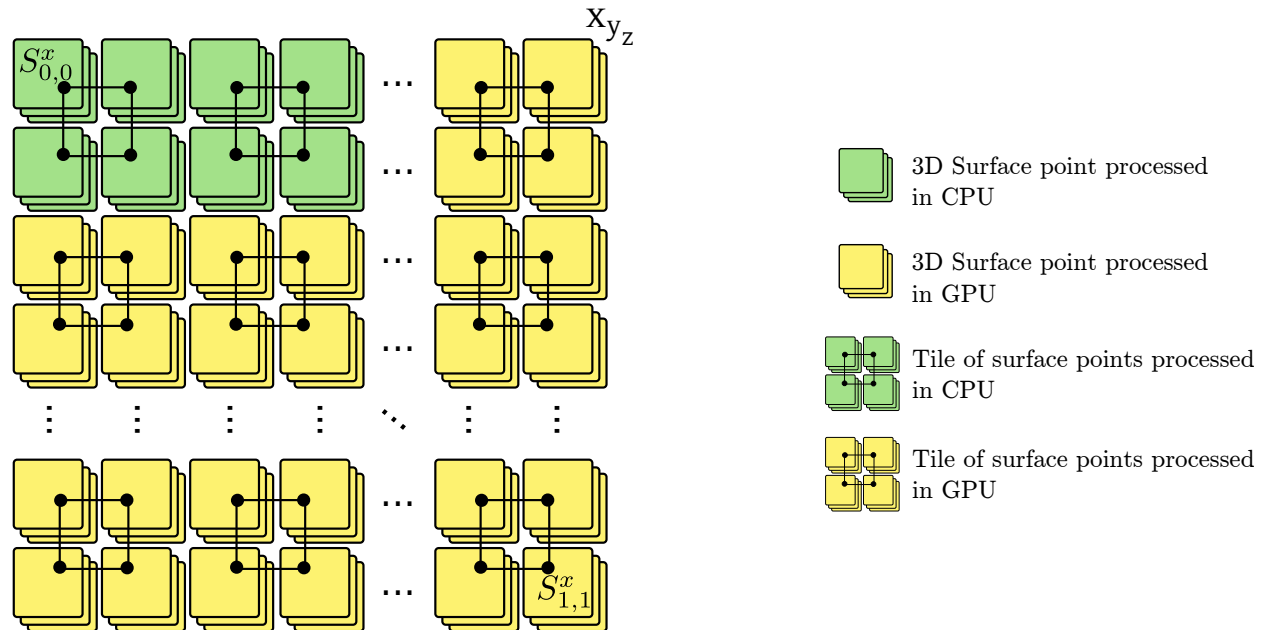
- Parametric non-rational formulation
 - Bernstein polynomials
 - Bi-cubic surface $m = n = 3$

$$\mathbf{S}(u, v) = \sum_{i=0}^m \sum_{j=0}^n \mathbf{P}_{i,j} B_{i,m}(u) B_{j,n}(v), \quad (1)$$

$$B_{i,m}(u) = \binom{m}{i} (1-u)^{(m-i)} u^i, \quad (2)$$

Bézier Surfaces: Static Distribution (I)

- Collaborative implementation
 - Tiles calculated by GPU blocks or CPU threads
 - **Static distribution**



Bézier Surfaces: Static Distribution (II)

■ Without Unified Memory

```
// Allocate control points
malloc(control_points, ...);
generate_cp(control_points);
cudaMalloc(d_control_points, ...);
cudaMemcpy(d_control_points, control_points, ..., HostToDevice); // Copy to device memory

// Allocate surface
malloc(surface, ...);
cudaMalloc(d_surface, ...);

// Launch CPU threads
std::thread main_thread (run_cpu_threads, control_points, surface, ...);

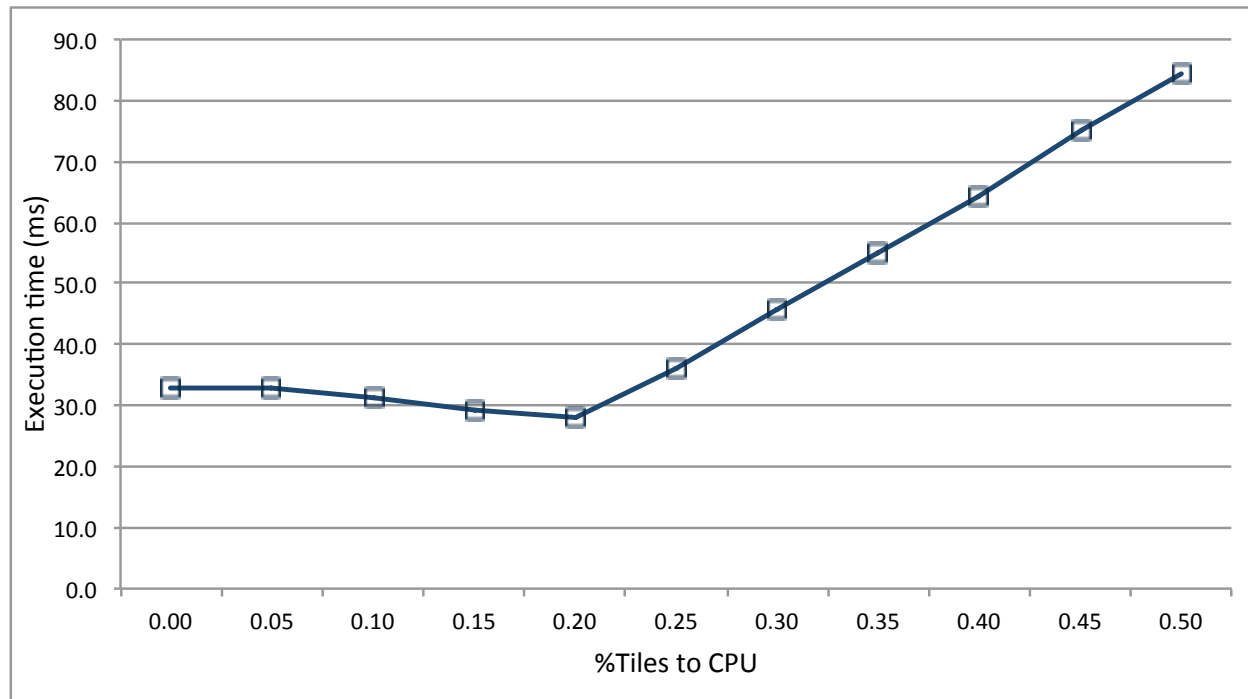
// Launch GPU kernel
gpu_kernel<<<blocks, threads>>> (d_surface, d_control_points, ...);

// Synchronize
main_thread.join();
cudaDeviceSynchronize();

// Copy GPU part of surface to host memory
cudaMemcpy(&surface[end_of_cpu_part], d_surface, ..., DeviceToHost);
```

Bézier Surfaces: Static Distribution (III)

- Performance results on NVIDIA Jetson TX1 (4 ARMv8 CPU cores + 2 GPU cores)
 - ❑ Bezier surface: 300x300, 4x4 control points
 - ❑ %Tiles to CPU
 - ❑ 17% speedup over GPU only



Bézier Surfaces with Unified Memory

■ With Unified Memory

```
// Allocate control points
malloc(control_points, ...);
generate_cp(control_points);
cudaMalloc(d_control_points, ...);
cudaMemcpy(d_control_points, control_points, ..., HostToDevice); // Copy to device memory

// Allocate surface
cudaMallocManaged(surface, ...);

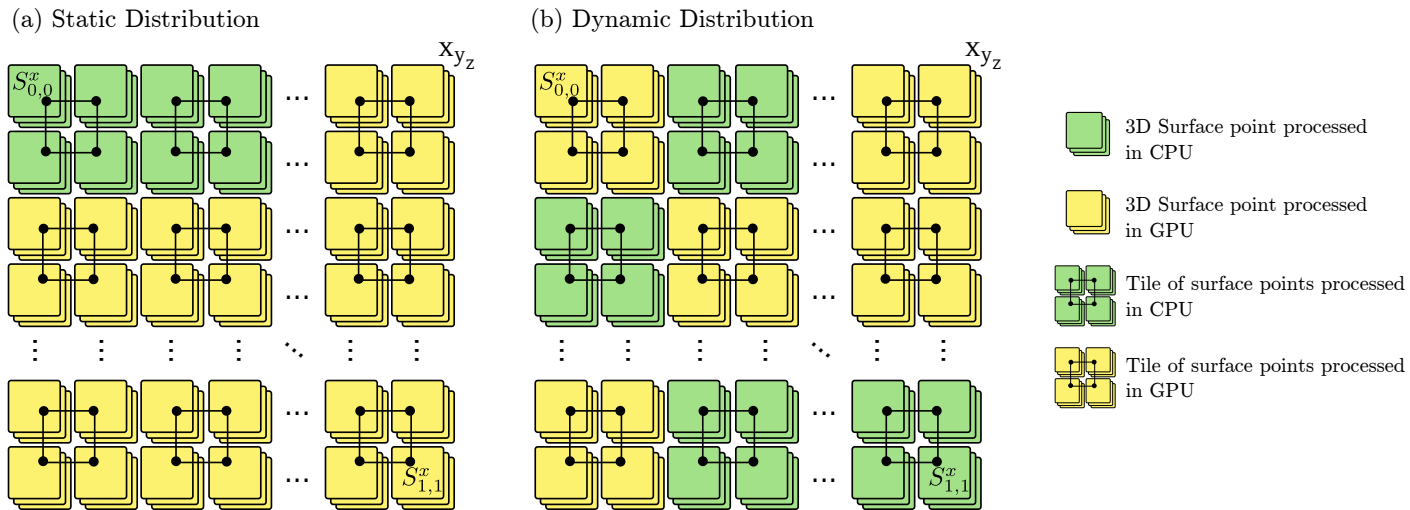
// Launch CPU threads
std::thread main_thread (run_cpu_threads, control_points, surface, ...);

// Launch GPU kernel
gpu_kernel<<<blocks, threads>>> (surface, d_control_points, ...);

// Synchronize
main_thread.join();
cudaDeviceSynchronize();
```


Bézier Surfaces: Dynamic Distribution

■ Static vs. dynamic implementation

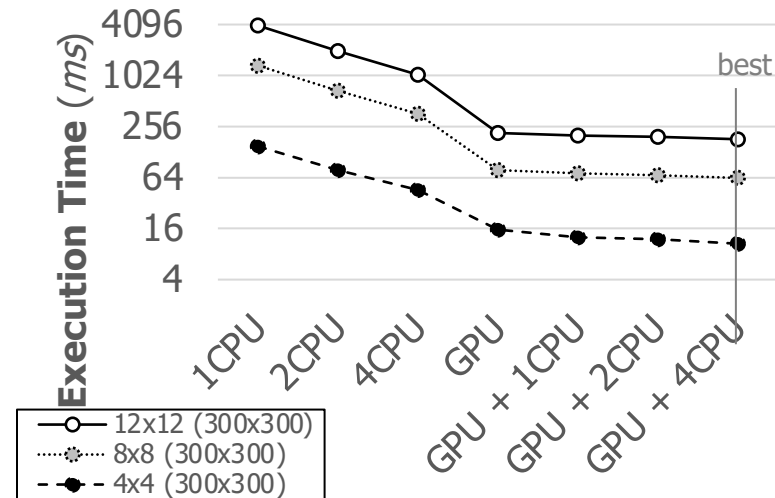


□ Pascal/Volta/Turing/Ampere Unified Memory: system-wide atomic operations

```
while(true){  
    if(threadIdx.x == 0)  
        my_tile = atomicAdd_system(tile_num, 1); // my_tile in shared memory; tile_num in UM  
  
    __syncthreads(); // Synchronization  
  
    if(my_tile >= number_of_tiles) break; // Break when all tiles processed  
  
    ... // Kernel body  
}
```

Benefits of Collaboration: Bézier Surfaces

- AMD Kaveri (4 CPU cores + 8 GPU cores)
 - Data partitioning improves performance



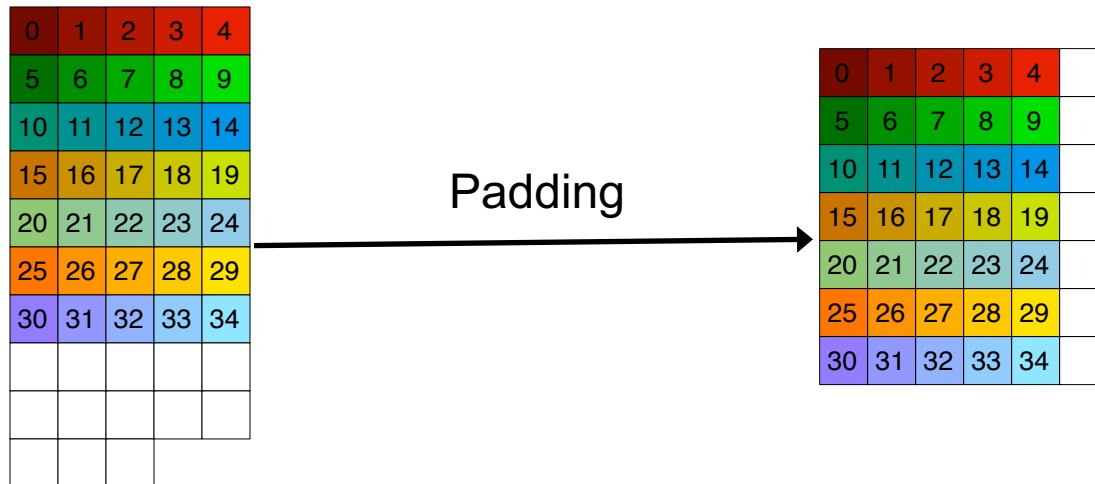
Bézier Surfaces
(up to 47% improvement over GPU only)

Padding (I)

- Matrix padding

- Use cases:

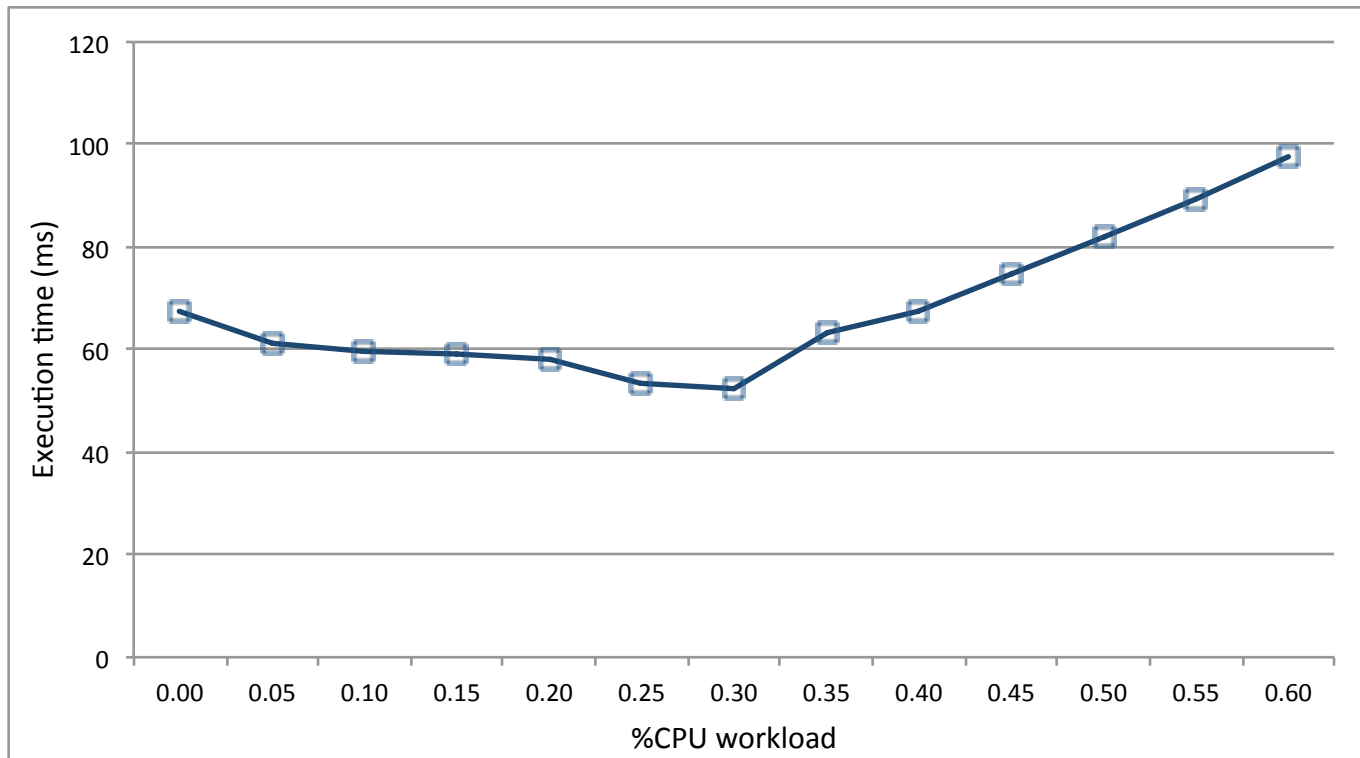
- Memory alignment
 - Transposition of near-square matrices
 - Etc.



- Traditionally, it can only be performed out-of-place

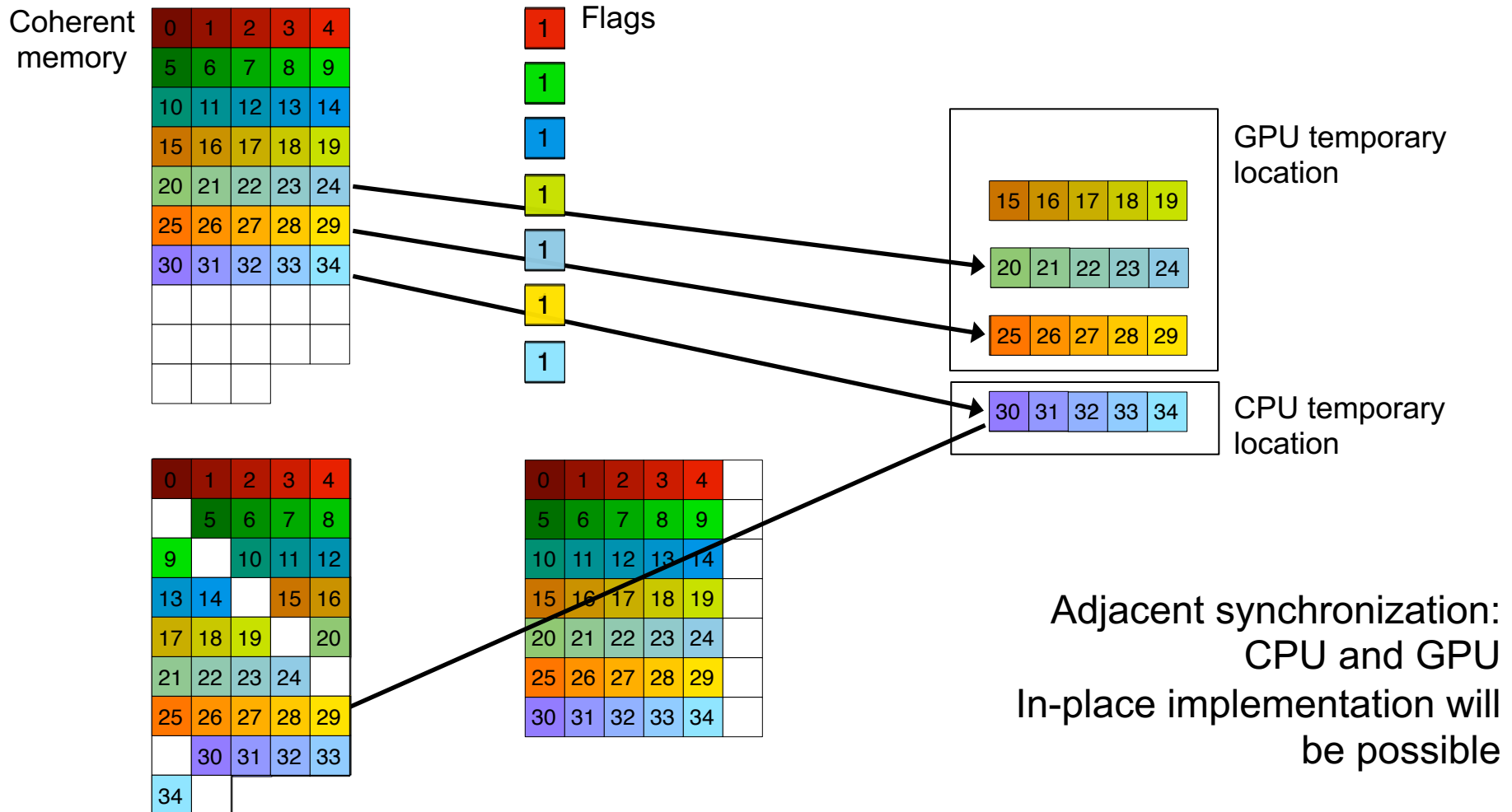
Padding (II)

- Performance results on NVIDIA Jetson TX1 (4 ARMv8 CPU cores + 2 GPU cores)
 - Matrix size: 4000x4000, padding = 1
 - 29% speedup over GPU only



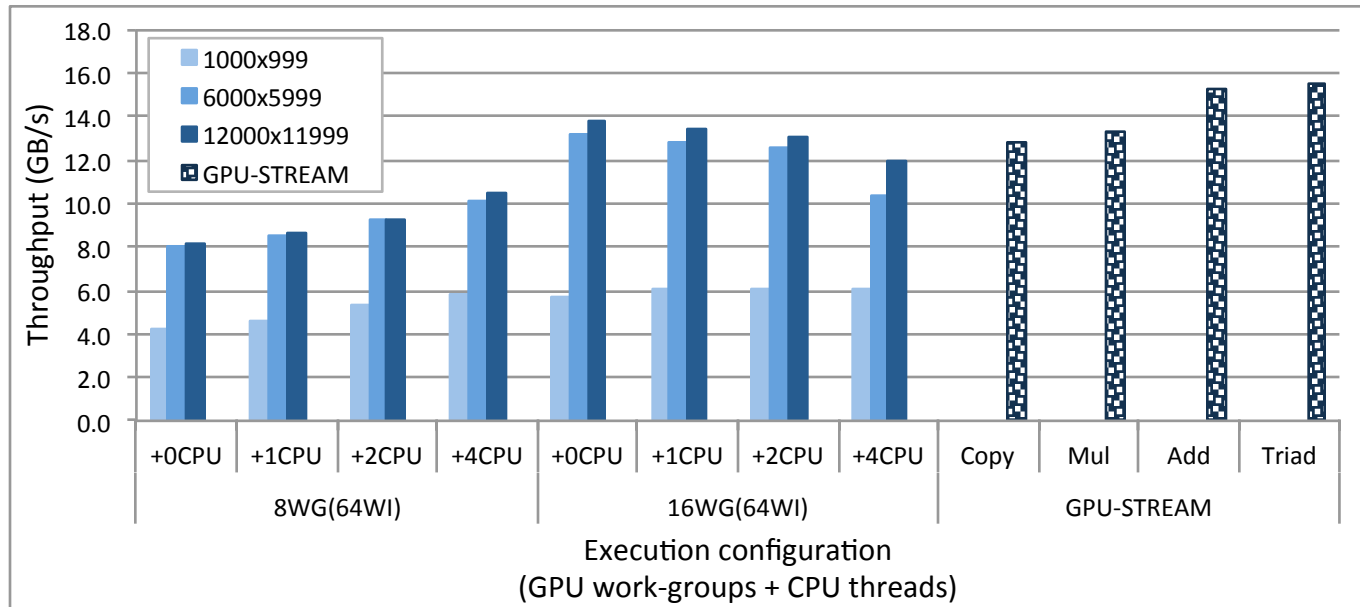
In-Place Padding

- Pascal/Volta/Turing/Ampere Unified Memory



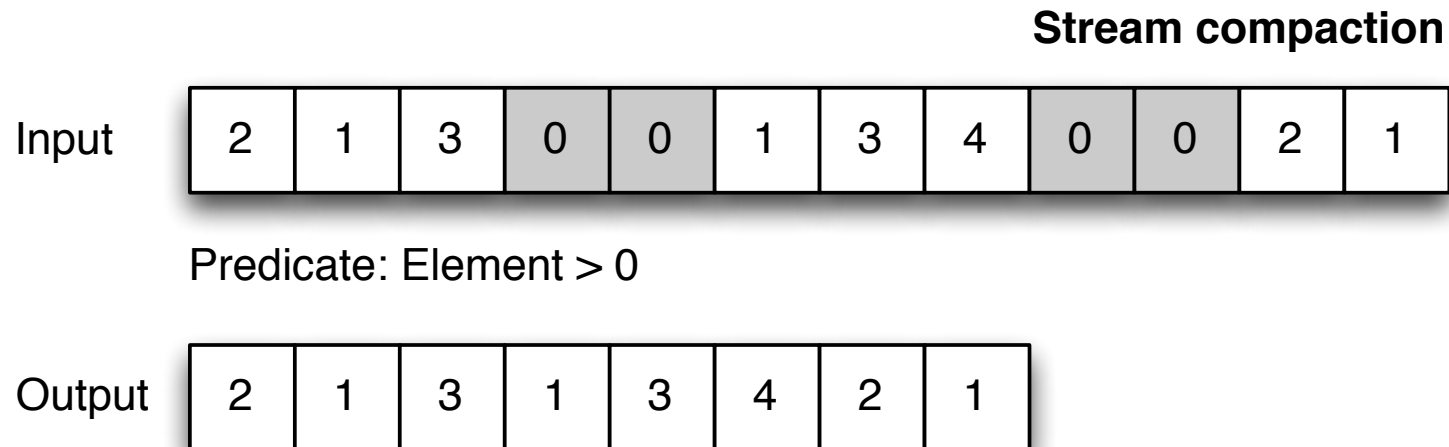
Benefits of Collaboration: Padding

- AMD Kaveri (4 CPU cores + 8 GPU cores)
 - Optimal number of devices is not always the maximum



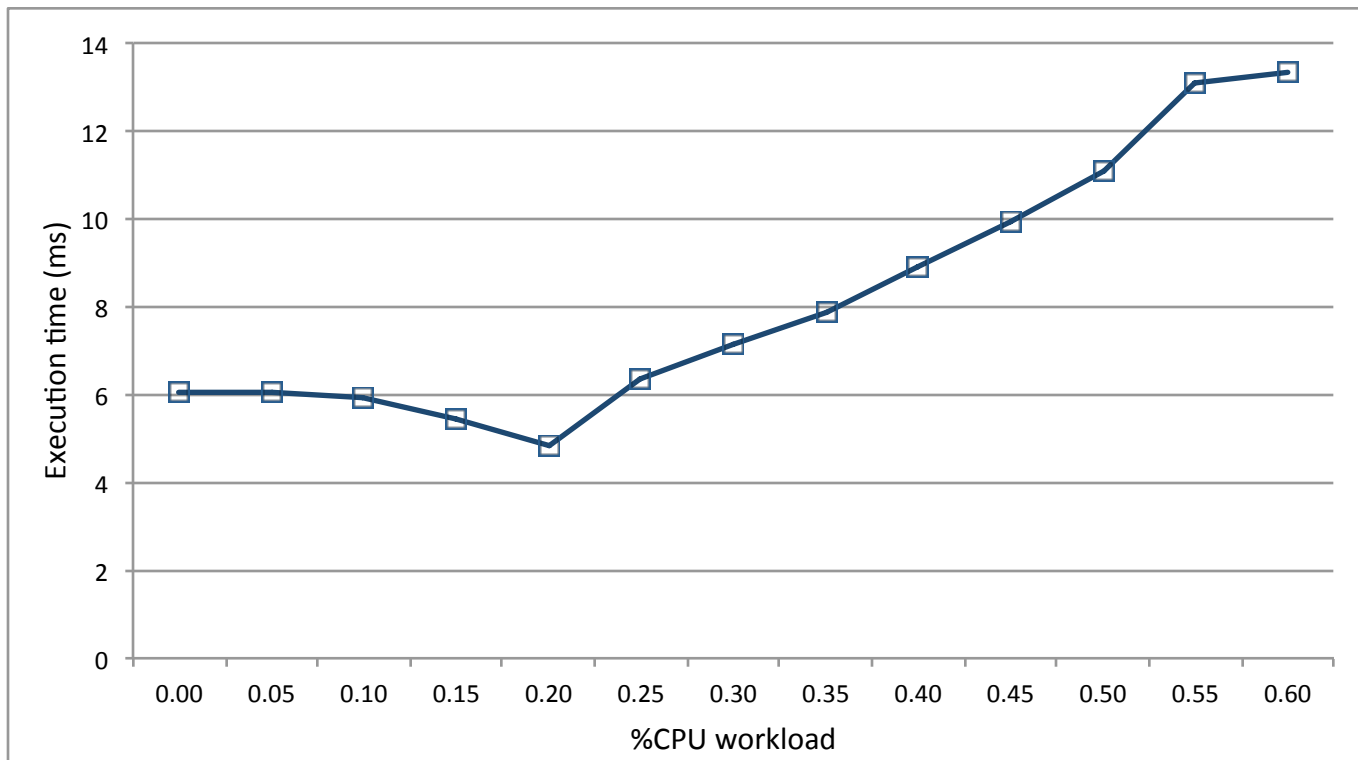
Stream Compaction (I)

- Stream compaction or filtering
 - Saving memory storage in sparse data
 - Similar to padding, but local reduction result (non-zero element count) is propagated



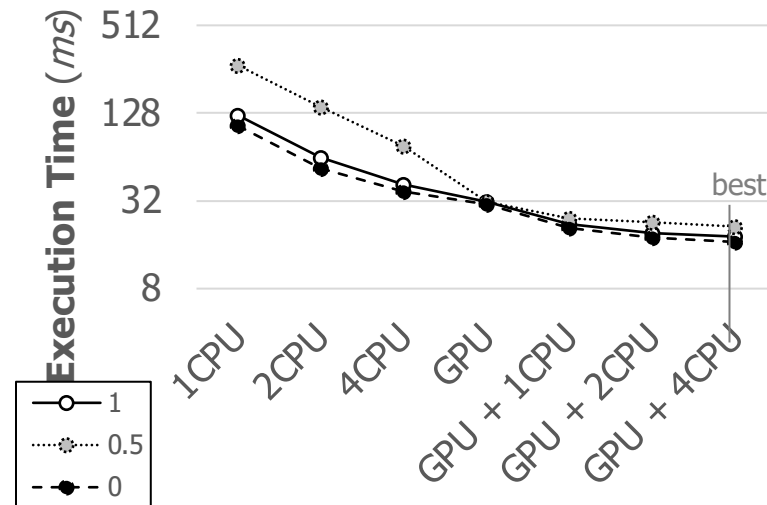
Stream Compaction (II)

- Performance results on NVIDIA Jetson TX1 (4 ARMv8 CPU cores + 2 GPU cores)
 - Array size: 2 MB, filtered items = 50%
 - 25% speedup over GPU only



Benefits of Collaboration: Stream Comp.

- AMD Kaveri (4 CPU cores + 8 GPU cores)
 - Data partitioning improves performance



Coarse-Grained Task Partitioning

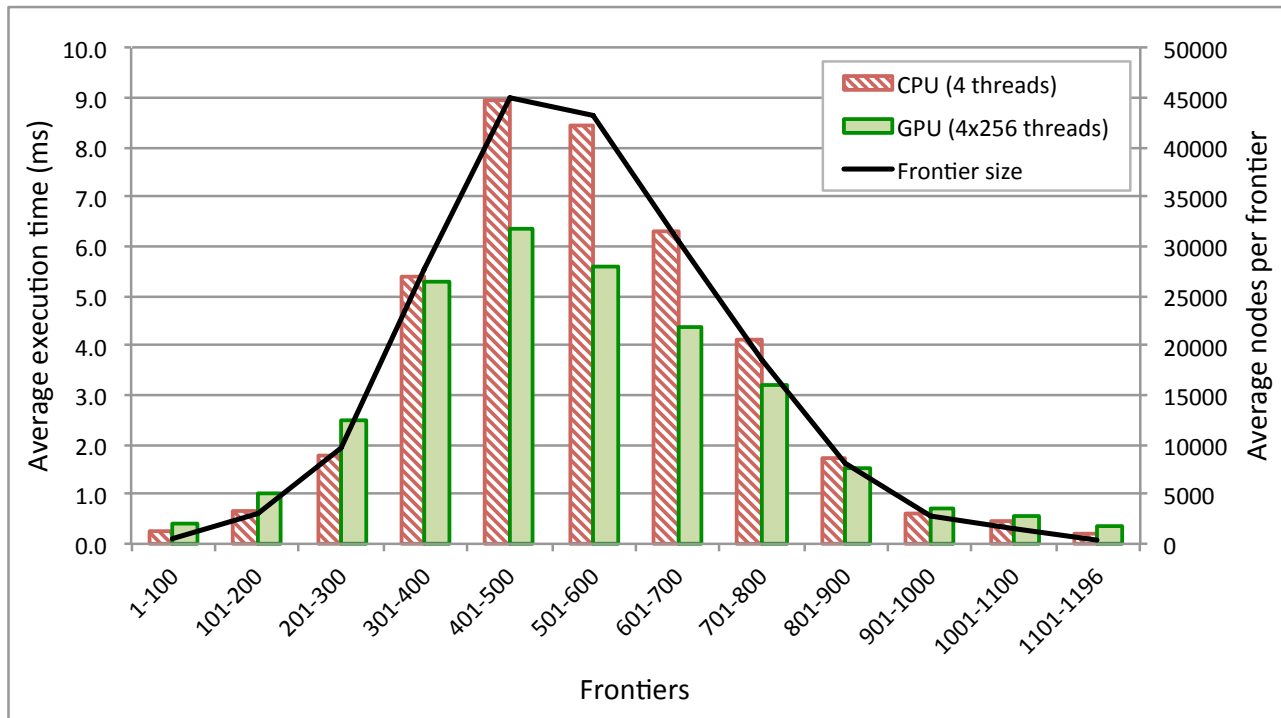
Breadth-First Search

- Small-sized and big-sized frontiers
 - Top-down approach
 - Kernel 1 and Kernel 2
- Atomic-based inter-block synchronization
 - Avoids kernel re-launch
- Very small frontiers
 - Underutilize GPU resources
- Collaborative implementation

Recall: BFS on CPU or GPU?

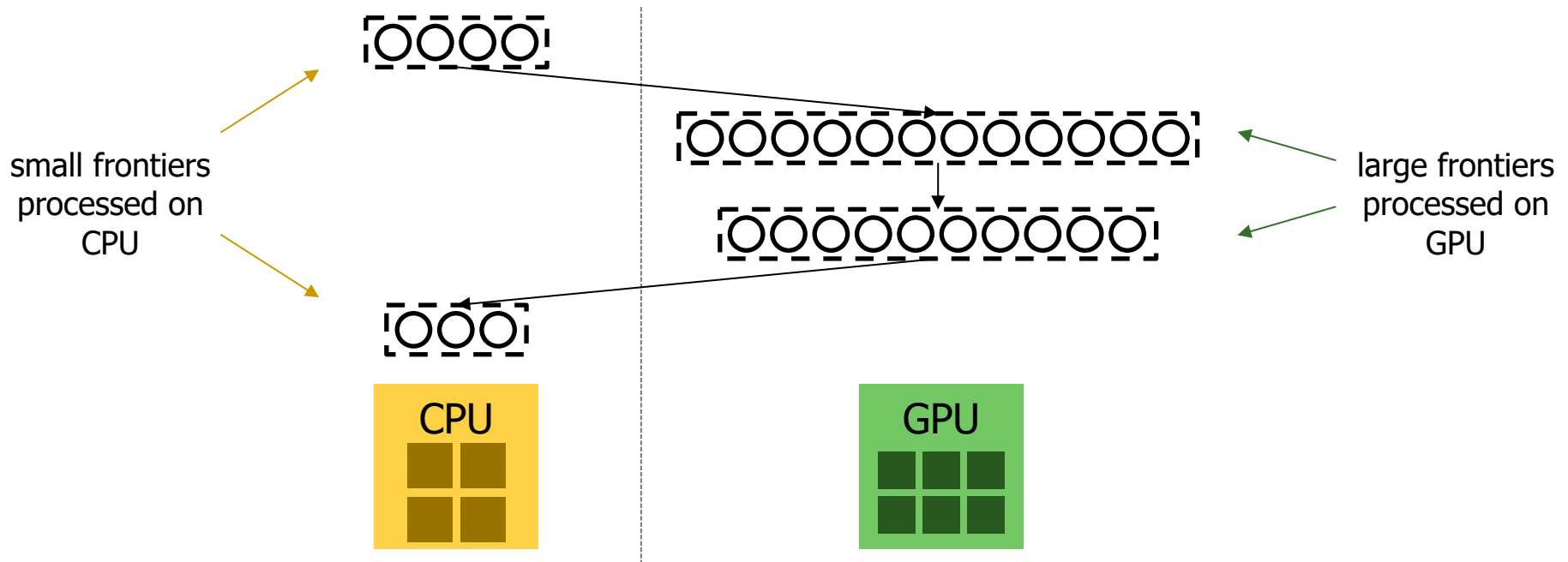
■ Motivation

- Small-sized frontiers underutilize GPU resources
 - NVIDIA Jetson TX1 (4 ARMv8 CPUs + 2 SMXs)
 - New York City roads



BFS: Collaborative Implementation

- Choose the most appropriate device



Collaborative Implementation without UM

- **Without** Unified Memory (UM)
 - Explicit memory copies

```
// Host code
while(frontier_size != 0){

    if(frontier_size < LIMIT){

        // Launch CPU threads

    }
    else{

        // Copy from host to device (queues and synchronization variables)

        // Launch GPU kernel

        // Copy from device to host (queues and synchronization variables)

    }

}
```

Collaborative Implementation with UM (I)

■ Unified Memory

- ❑ `cudaMallocManaged()`;
- ❑ Easier programming
- ❑ No explicit memory copies

```
// Host code
while(frontier_size != 0){

    if(frontier_size < LIMIT){

        // Launch CPU threads

    }
    else{

        // Launch GPU kernel for every frontier (kernel termination and relaunch)

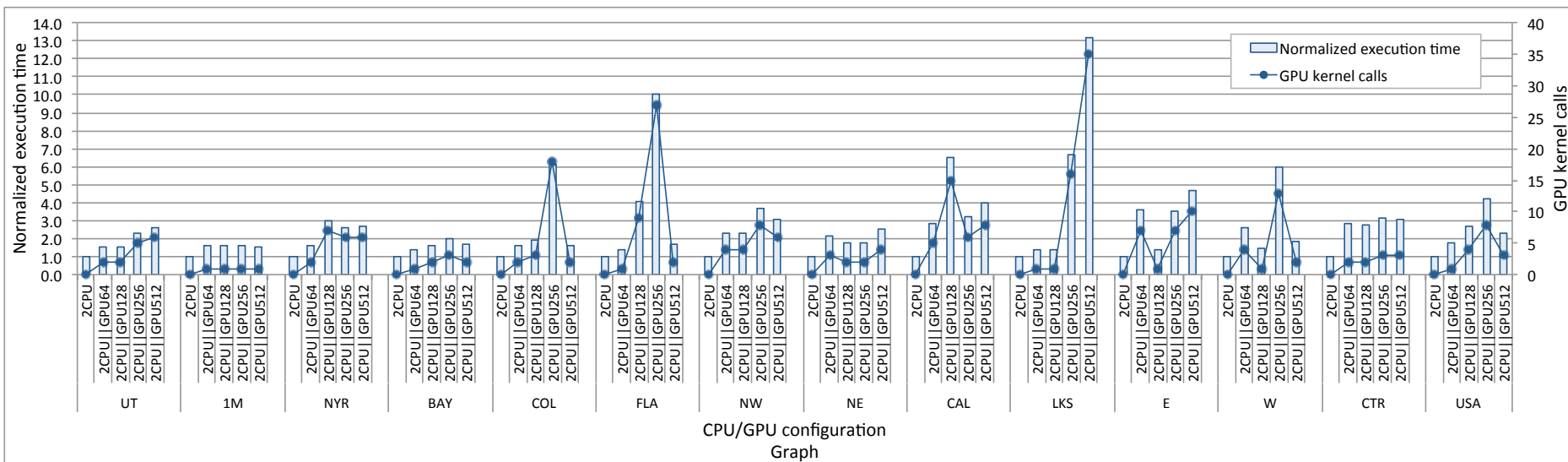
        cudaDeviceSynchronize();

    }

}
```

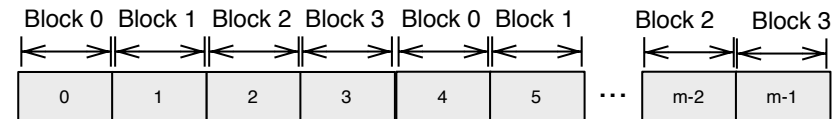
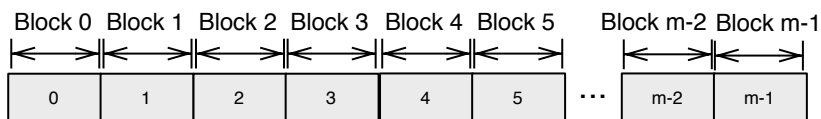
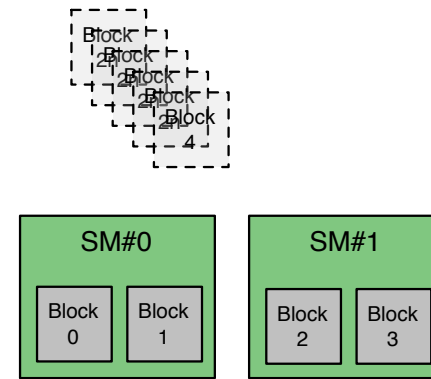
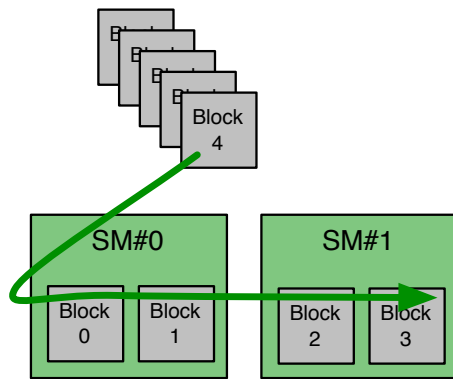
BFS: Kernel Termination and Relaunch

- AMD Kaveri (4 CPU cores + 8 GPU cores)
 - High overhead of kernel relaunch makes CPU+GPU collaboration impractical



Recall: Persistent Thread Blocks

- Combine Kernel 1 and Kernel 2
- We can **avoid kernel re-launch**
- We need to use **persistent thread blocks**
 - Kernel 2 launches ($\text{frontier_size} / \text{block_size}$) blocks
 - Persistent blocks: up to $(\text{number_SMs} \times \text{max_blocks_SM})$



Atomic-based Block Synchronization (I)

■ Code (simplified)

```
// GPU kernel
const int gtid = blockIdx.x * blockDim.x + threadIdx.x;

while(frontier_size != 0){

    for(node = gtid; node < frontier_size; node += blockDim.x * gridDim.x){

        // Visit neighbors
        // Enqueue in output queue if needed (global or local queue)

    }

    // Update frontier_size

    // Global synchronization
}
```

Atomic-based Block Synchronization (II)

■ Global synchronization (simplified)

□ At the end of each iteration

```
const int tid = threadIdx.x;
const int gtid = blockIdx.x * blockDim.x + threadIdx.x;
atomicExch(ptr_threads_run, 0);
atomicExch(ptr_threads_end, 0);
int frontier = 0;
...
frontier++;

if(tid == 0){
    atomicAdd(ptr_threads_end, 1); // Thread block finishes iteration
}

if(gtid == 0){
    while(atomicAdd(ptr_threads_end, 0) != gridDim.x){;} // Wait until all blocks finish

    atomicExch(ptr_threads_end, 0); // Reset
    atomicAdd(ptr_threads_run, 1); // Count iteration
}

if(tid == 0 && gtid != 0){
    while(atomicAdd(ptr_threads_run, 0) < frontier){;} // Wait until ptr_threads_run is updated
}

__syncthreads(); // Rest of threads wait here

...
```

BFS: Collaborative Implementation (II)

- Choose CPU or GPU depending on frontier

```
// Host code
while(frontier_size != 0){

    if(frontier_size < LIMIT){

        // Launch CPU threads
    }
    else{

        // Launch GPU kernel (keep running while frontier_size >= LIMIT)

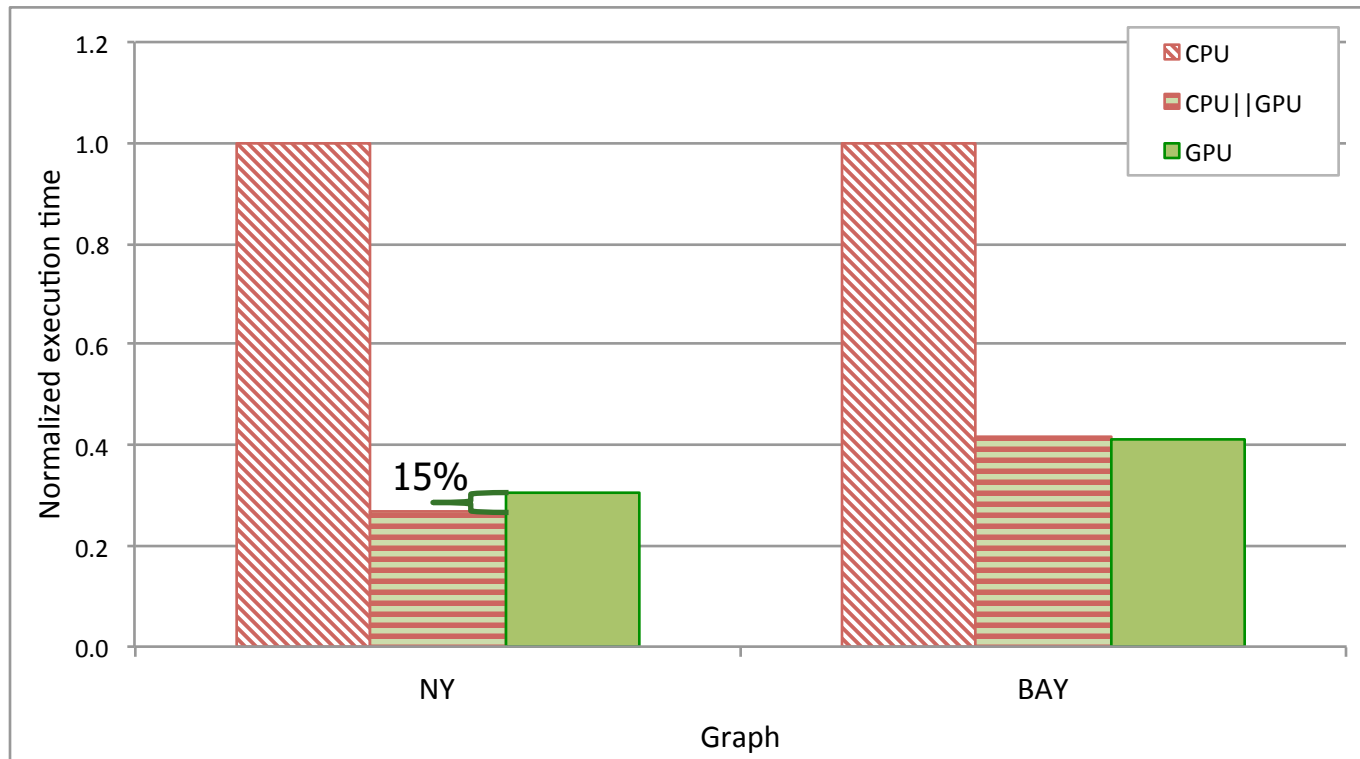
        cudaDeviceSynchronize();
    }
}
```

- CPU threads or GPU kernel keep running while the condition is satisfied

BFS: Collaborative Implementation (III)

■ Experimental results

- NVIDIA Jetson TX1 (4 ARMv8 CPU cores + 2 GPU cores)



Collaborative Implementation with UM (II)

- Pascal/Volta/Turing/Ampere Unified Memory & HSA
 - ❑ CPU/GPU coherence
 - ❑ System-wide atomic operations
 - ❑ No need to re-launch kernel or CPU threads
 - ❑ Possibility of CPU and GPU working on the same frontier

```
// Host code
while(frontier_size != 0){

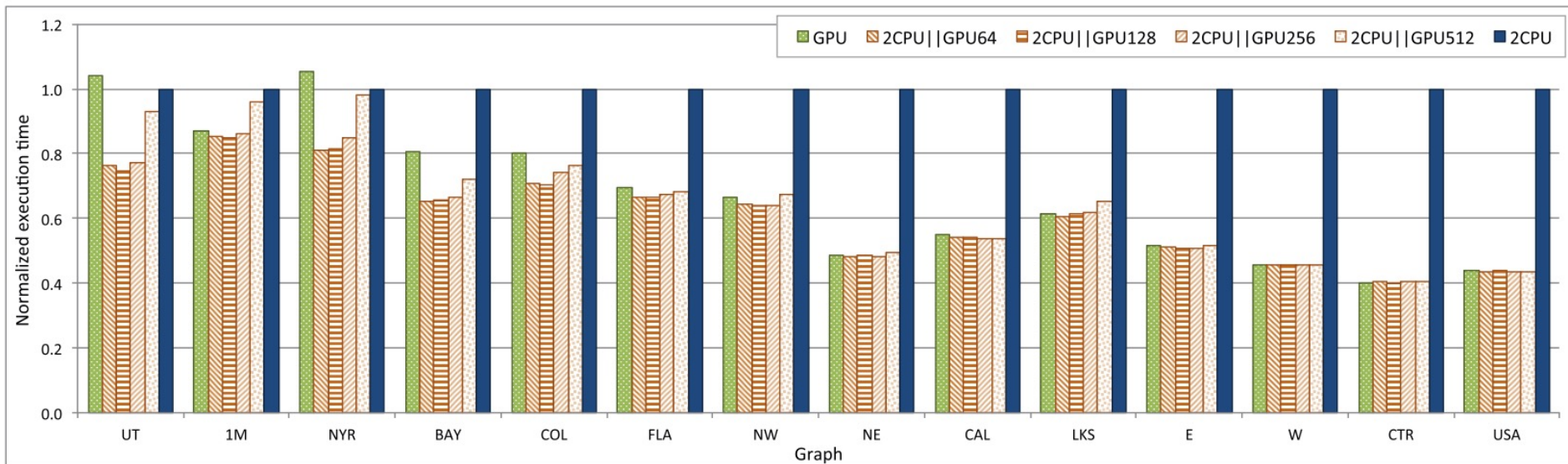
    if(frontier_size < LIMIT){

        // Launch CPU threads (compute when frontier_size < LIMIT)
    }
    else{

        // Launch GPU kernel (compute when frontier_size >= LIMIT)
    }
}
cudaDeviceSynchronize();
```

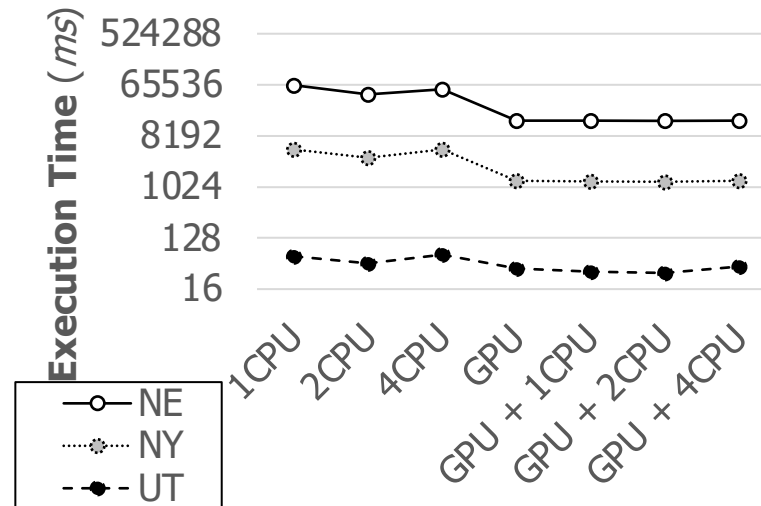
Benefits of Collaboration: BFS

- AMD Kaveri (4 CPU cores + 8 GPU cores)
 - The collaborative implementation (with system-wide atomics) is up to 39% faster than the GPU only version



Benefits of Collaboration: SSSP

- AMD Kaveri (4 CPU cores + 8 GPU cores)
 - SSSP performs more computation than BFS



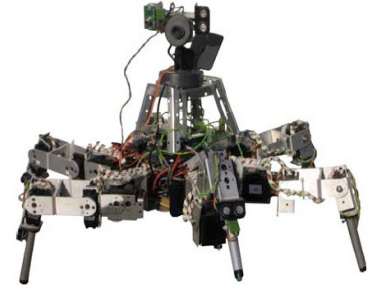
Single Source Shortest Path
(up to 22% improvement over GPU only)

Fine-Grained Task Partitioning

Egomotion Compensation and Moving Objects Detection (I)

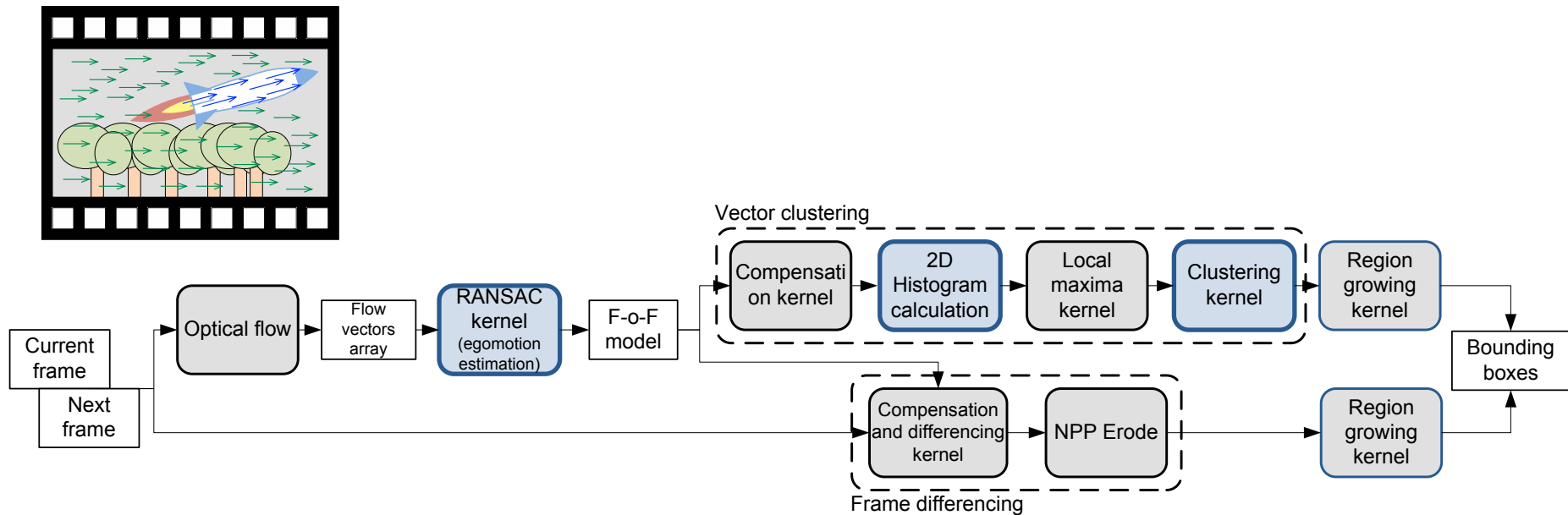
■ Hexapod robot OSCAR

- Rescue scenarios
- Strong egomotion on uneven terrains



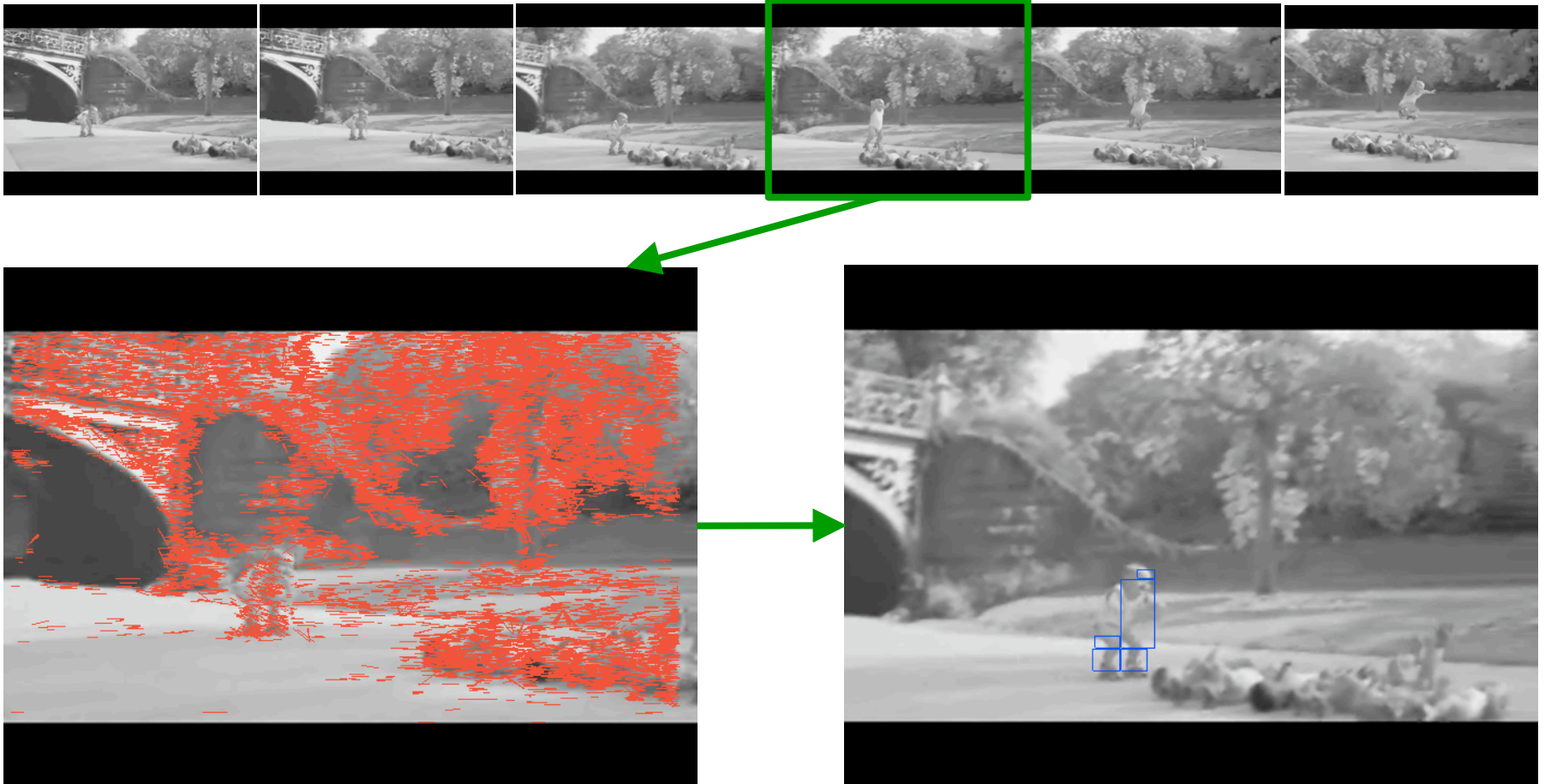
■ Algorithm

- Random Sample Consensus (RANSAC): F-o-F model



Egomotion Compensation and Moving Objects Detection (II)

Fast moving object in strong egomotion scenario detected by vector clustering

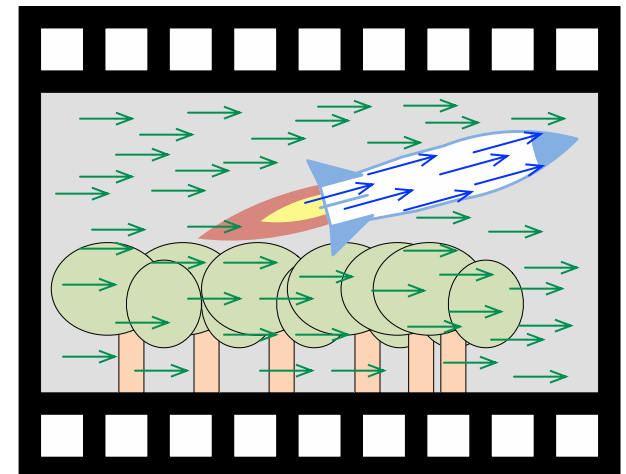


RANSAC: SISD and SIMD Phases

■ RANSAC (Fischler+, 1981)

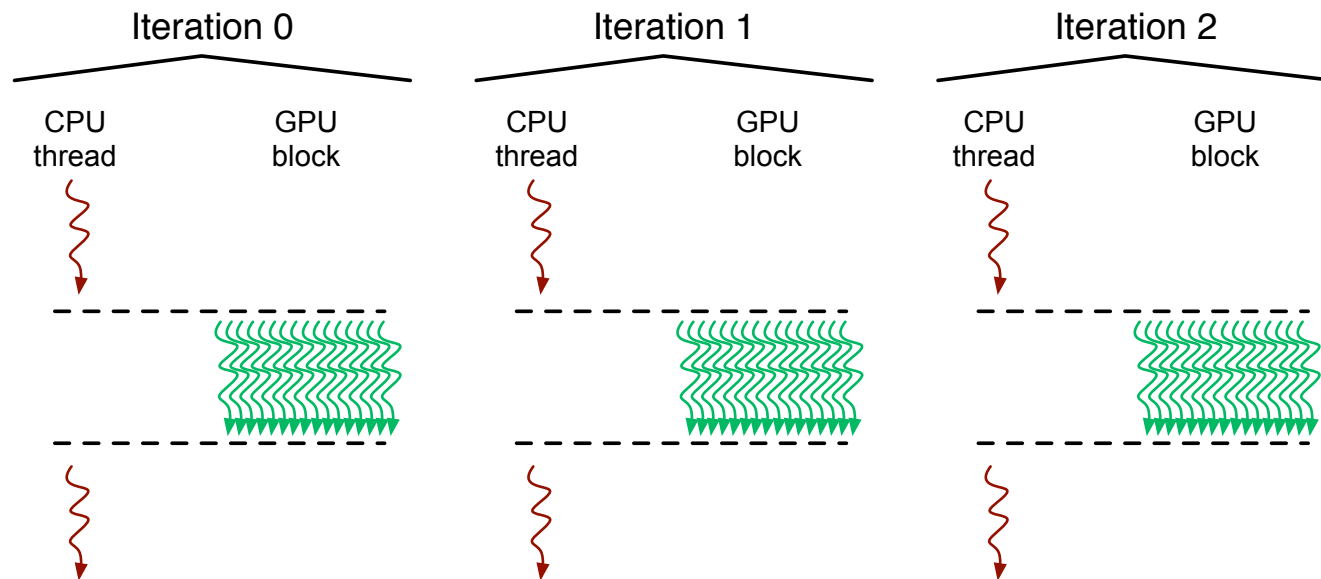
```
while (iteration < MAX_ITER){  
    Fitting stage (Compute F-o-F model)           // SISD phase  
    Evaluation stage (Count outliers)             // SIMD phase  
    Comparison to best model                      // SISD phase  
    Check if best model is good enough and iteration >= MIN_ITER // SISD phase  
}
```

- ❑ Fitting stage picks two flow vectors randomly
- ❑ Evaluation generates motion vectors from F-o-F model, and compares them to real flow vectors

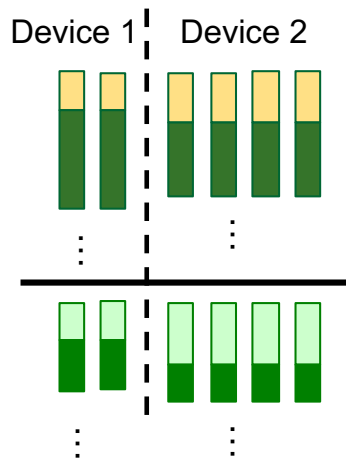
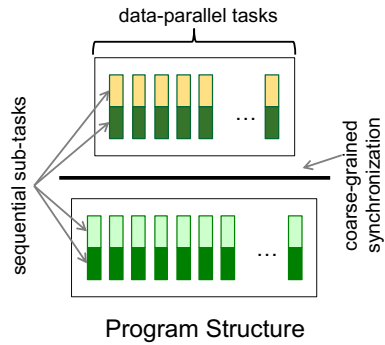


Collaborative Implementation

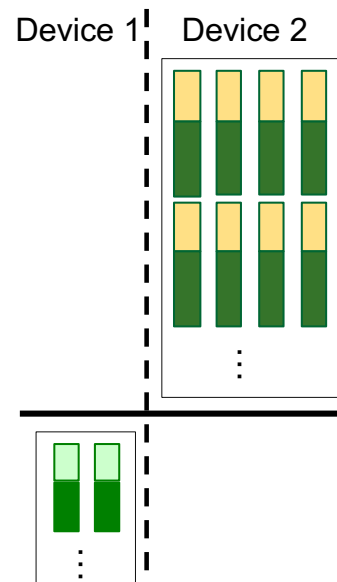
- Randomly picked vectors: Iterations are independent
 - We assign one iteration to one CPU thread and one GPU block



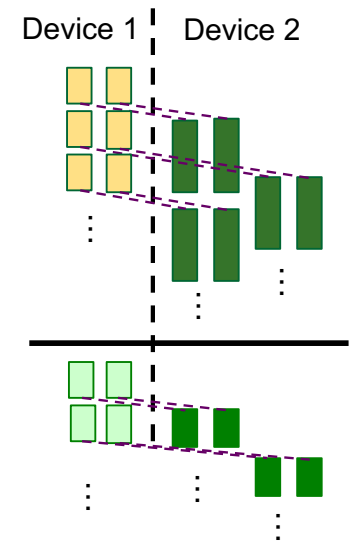
Collaborative Patterns



Data Partitioning



Coarse-grained Task Partitioning



Fine-grained Task Partitioning

Chai Benchmark Suite

- Collaborative Heterogeneous Applications for Integrated architectures
- Heterogeneous execution on CPU, GPU, FPGA
- Collaboration patterns
 - 8 data partitioning benchmarks
 - 3 coarse-grain task partitioning benchmarks
 - 3 fine-grain task partitioning benchmarks
- Discrete (D) and Unified (U) versions
 - CUDA, OpenCL, and C++AMP for CPU+GPU
 - OpenCL for CPU+FPGA
 - CUDA-Sim for Gem5-GPU



<https://chai-benchmarks.github.io>



Chai Benchmarks

Collaboration Pattern		Short Name	Benchmark
Data Partitioning		BS	Bézier Surface
		CEDD	Canny Edge Detection
		HSTI	Image Histogram (Input Partitioning)
		HSTO	Image Histogram (Output Partitioning)
		PAD	Padding
		RSCD	Random Sample Consensus
		SC	Stream Compaction
		TRNS	In-place Transposition
Task Partitioning	Fine-grain	RSCT	Random Sample Consensus
		TQ	Task Queue System (Synthetic)
		TQH	Task Queue System (Histogram)
	Coarse-grain	BFS	Breadth-First Search
		CEDT	Canny Edge Detection
		SSSP	Single-Source Shortest Path

Versions:

- OpenCL-**U**
- OpenCL-**D**
- CUDA-**U**
- CUDA-**D**
- CUDA-**U**-Sim
- CUDA-**D**-Sim
- C++AMP

Chai: Diversity of Benchmarks (I)

- Diversity of partitioning, usage of system-wide atomics, load balancing, and concurrency

DATA PARTITIONING

Benchmark	Partitioning Granularity	Partitioned Data	System-wide Atomics	Load Balance
BS	Fine	Output	None	Yes
CEDD	Coarse	Input, Output	None	Yes
HSTI	Fine	Input	Compute	No
HSTO	Fine	Output	None	No
PAD	Fine	Input, Output	Sync	Yes
RSCD	Medium	Output	Compute	Yes
SC	Fine	Input, Output	Sync	No
TRNS	Medium	Input, Output	Sync	No

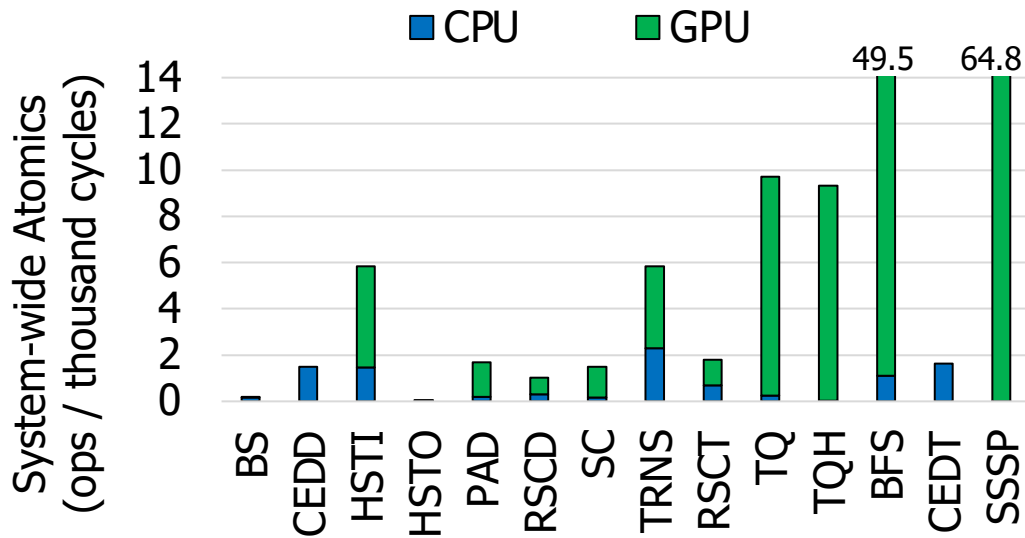
FINE-GRAIN TASK PARTITIONING

Benchmark	System-wide Atomics	Load Balance
RSCT	Sync, Compute	Yes
TQ	Sync	No
TQH	Sync	No

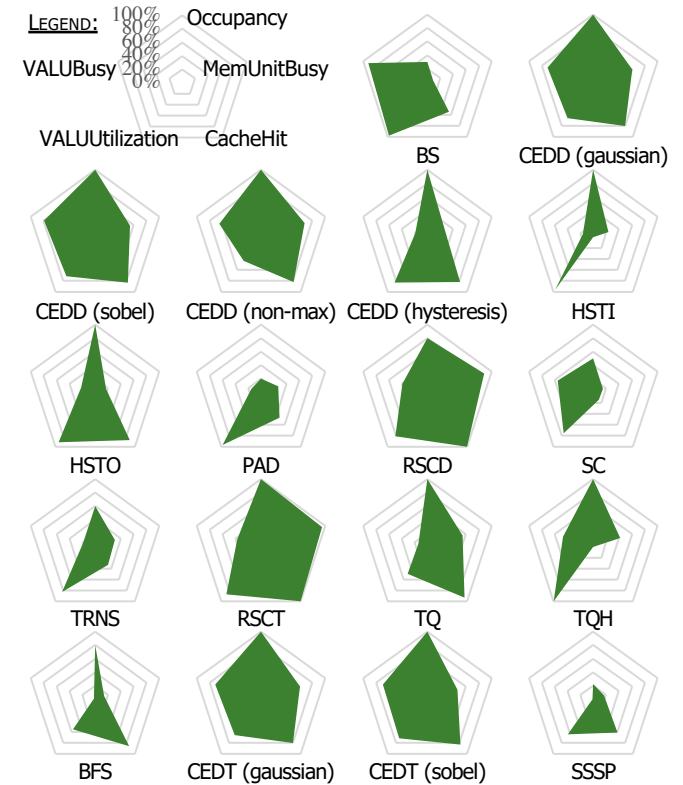
COARSE-GRAIN TASK PARTITIONING

Benchmark	System-wide Atomics	Partitioning	Concurrency
BFS	Sync, Compute	Iterative	No
CEDT	Sync	Non-iterative	Yes
SSSP	Sync, Compute	Iterative	No

Chai: Diversity of Benchmarks (II)

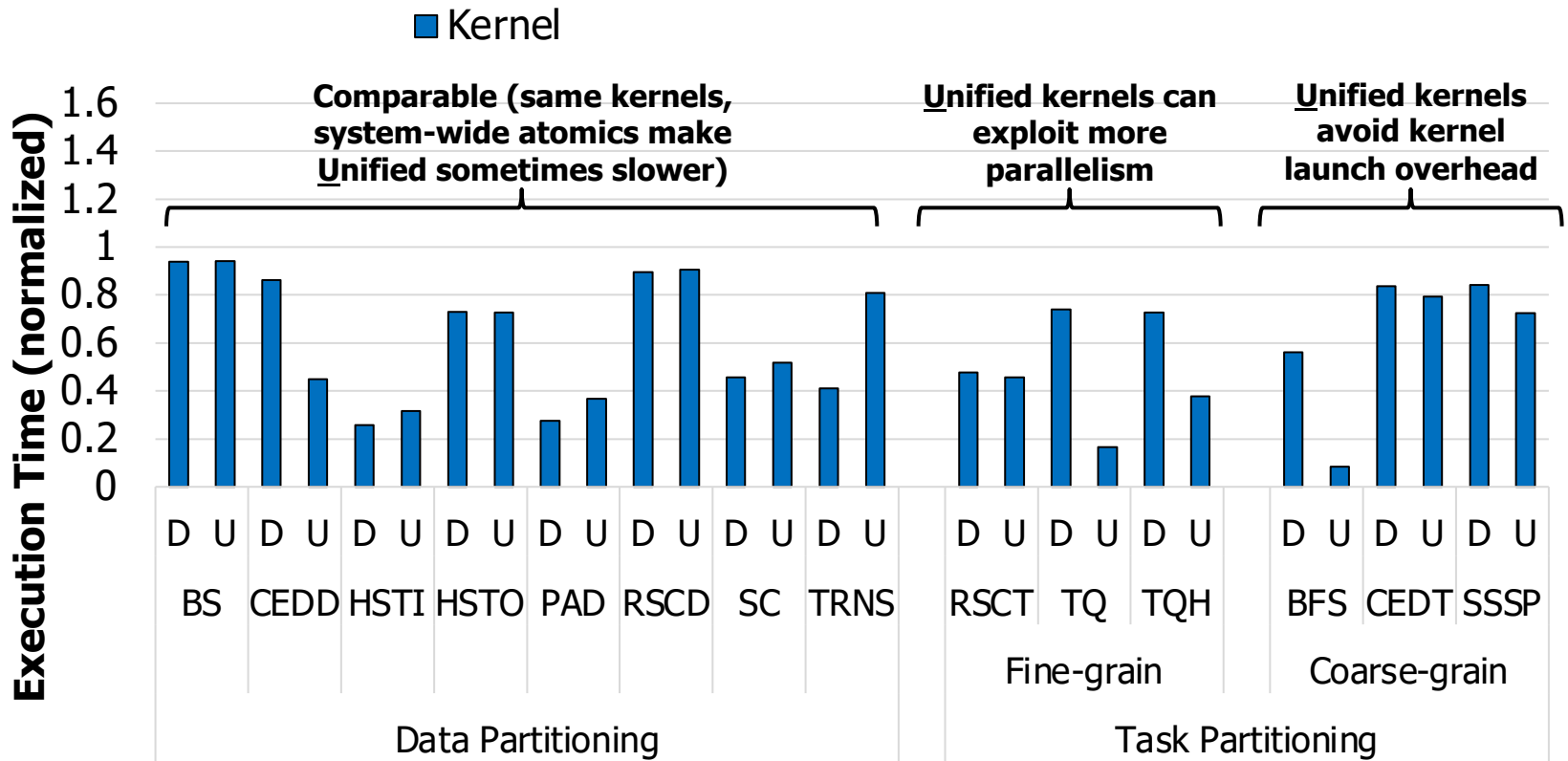


Varying intensity in use of system-wide atomics



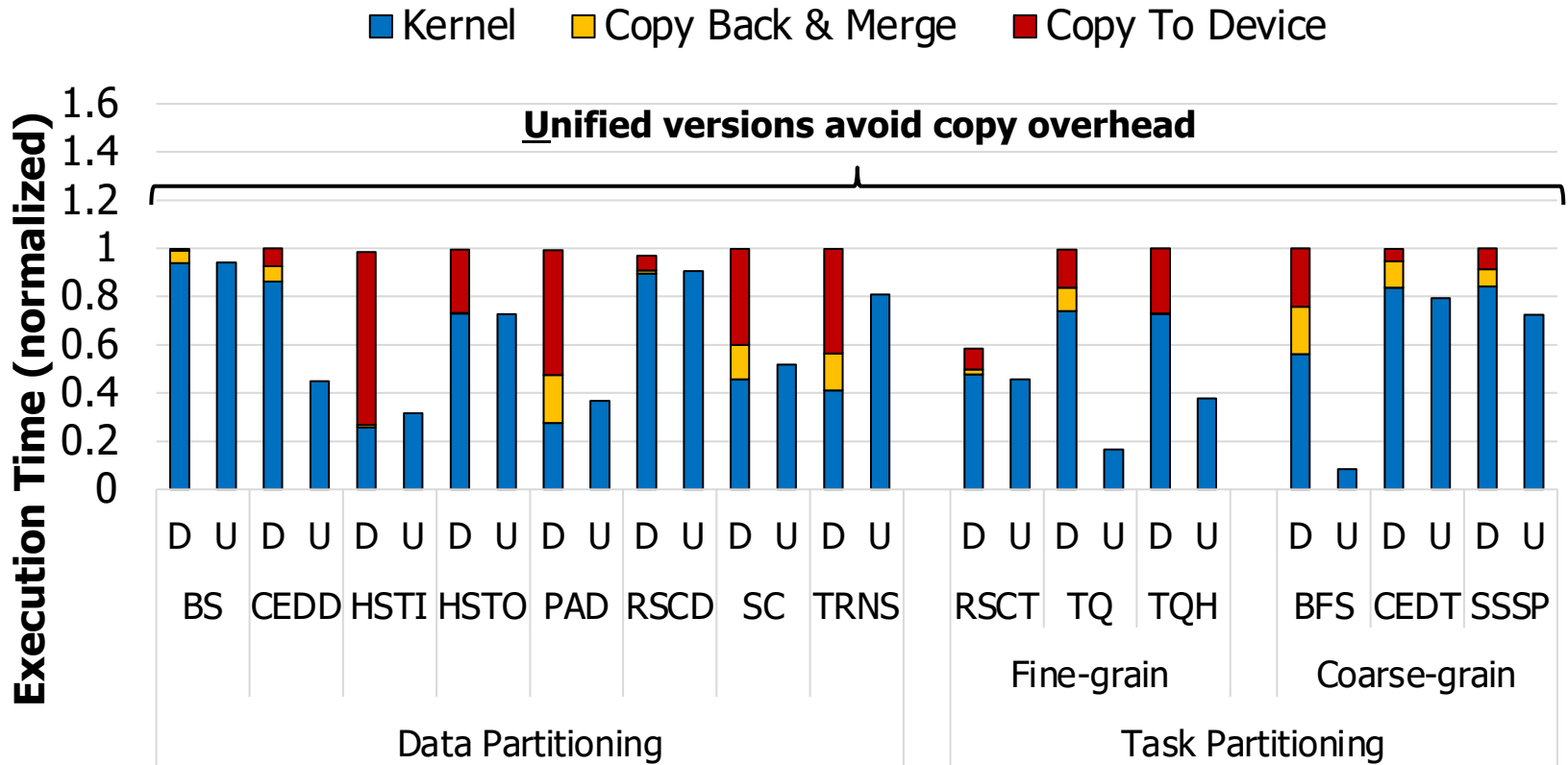
Diverse execution profiles

Benefits of Unified Memory: Kernel Time



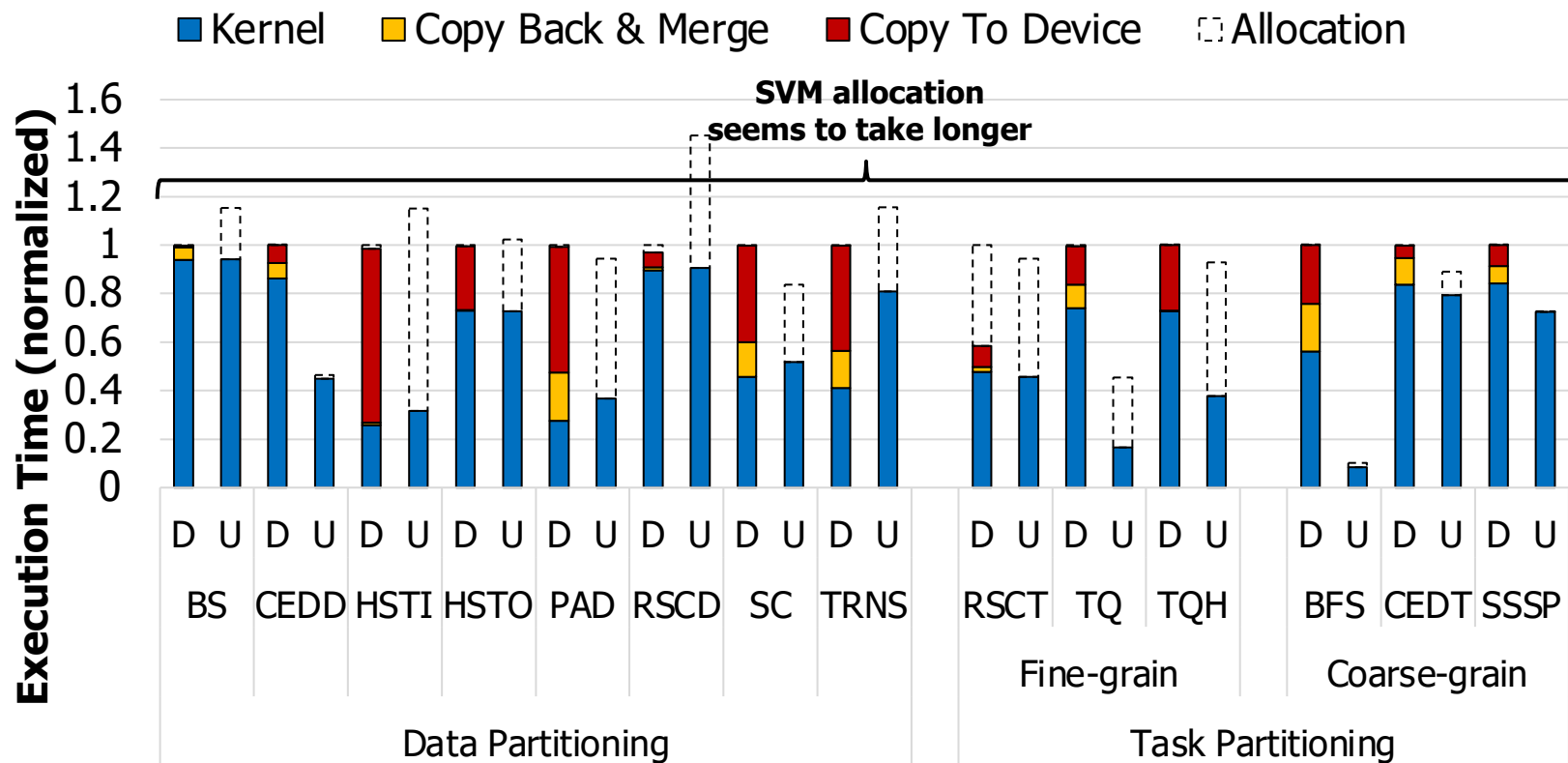
AMD Kaveri (4 CPU cores + 8 GPU cores), OpenCL

Benefits of Unified Memory: Data Transfers



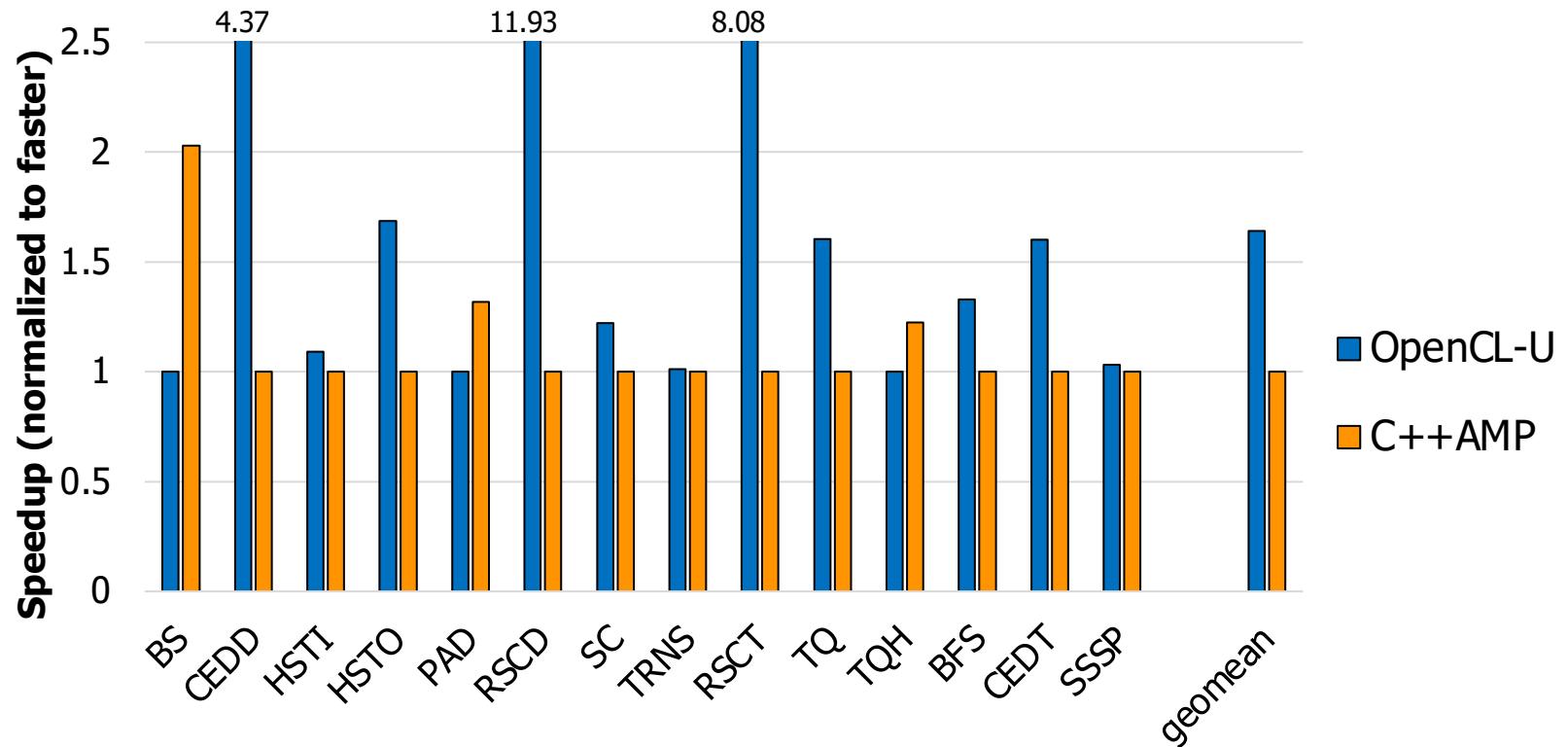
AMD Kaveri (4 CPU cores + 8 GPU cores), OpenCL

Benefits of Unified Memory: Allocation



AMD Kaveri (4 CPU cores + 8 GPU cores), OpenCL

Comparison C++AMP vs. OpenCL-U

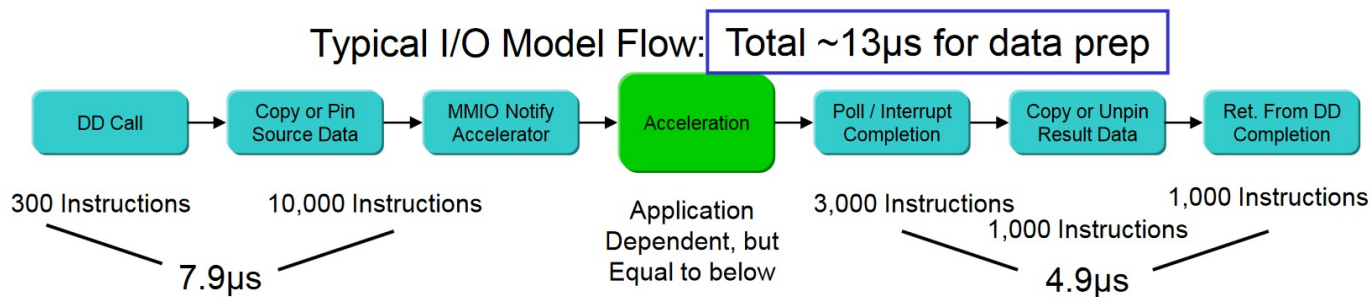
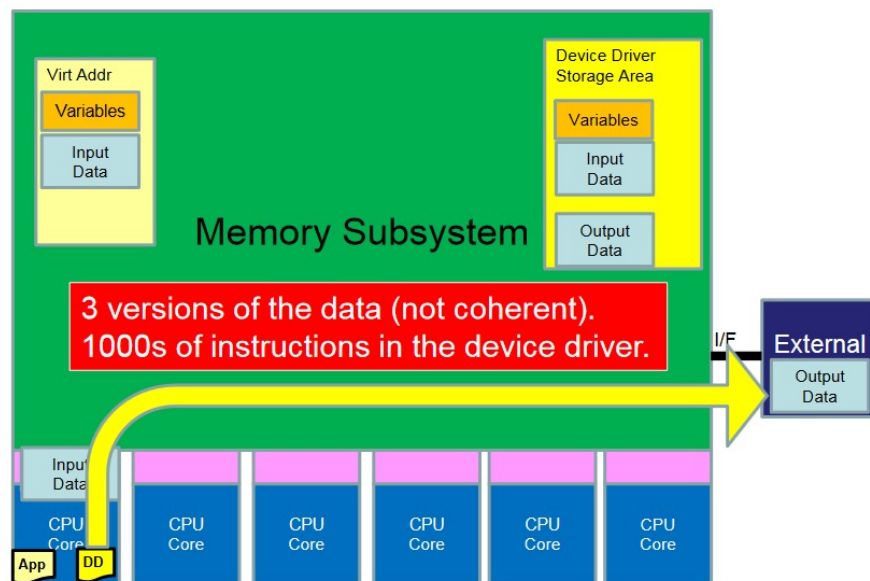


Heterogeneous System Architecture

- Wen-mei W. Hwu (editor), “**Heterogeneous System Architecture: A New Compute Platform Infrastructure,**” 2016
 - Chapter 8 – Application use cases: Platform atomics



Background: Traditional I/O Technology

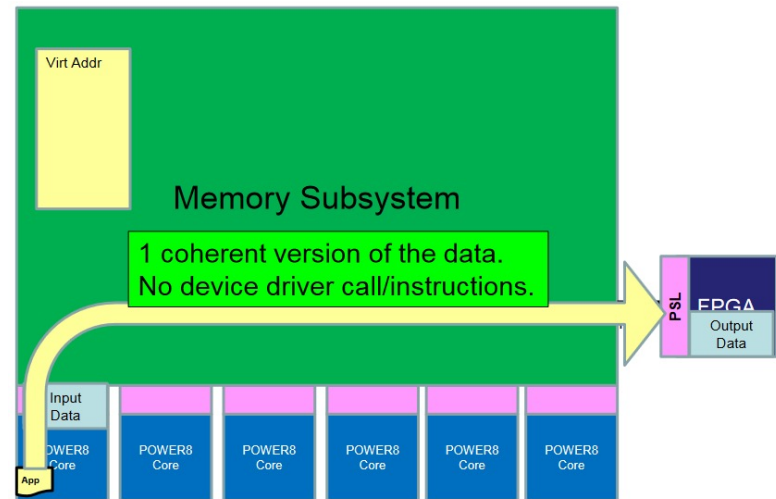
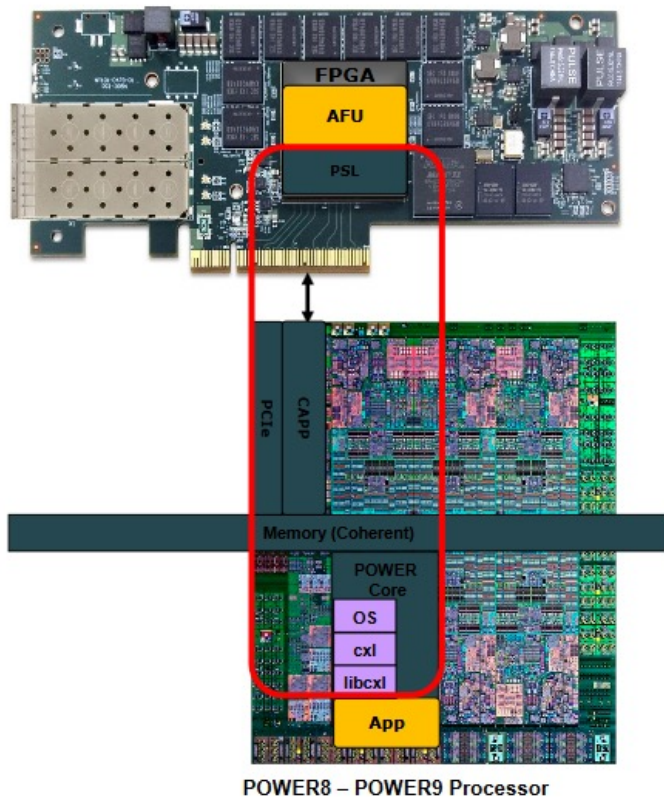


Dionysios Diamantopoulos, IBM Research – Zurich, COOL Chips 2018

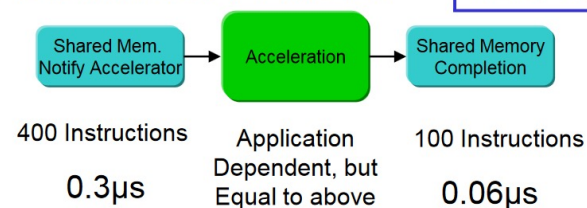


CAPI/OpenCAPI Overview

- CAPI/CAPI2 (Coherent Accelerator Processor Interface)
- OpenCAPI



Flow with a CAPI Model:



Total 0.36μs

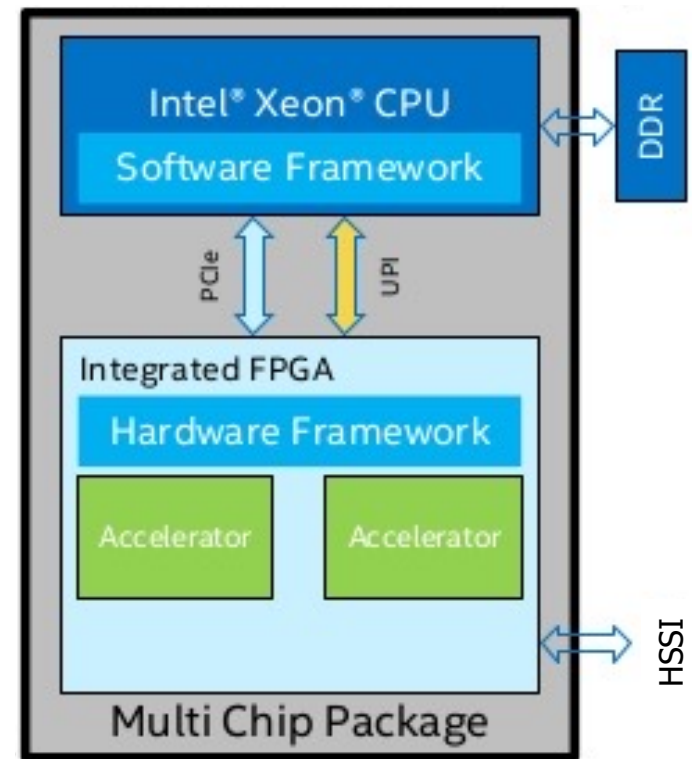
Dionysios Diamantopoulos, IBM Research – Zurich, COOL Chips 2018



Collaborative Computing on CPU+FPGA

- Traditionally, accelerators (GPUs, FPGAs, etc.) have been used as *offload* engines
- Heterogeneous architectures moving towards tighter integration
 - ❑ Unified memory
 - ❑ System-wide atomics
- Tighter integration allows fine-grained collaboration

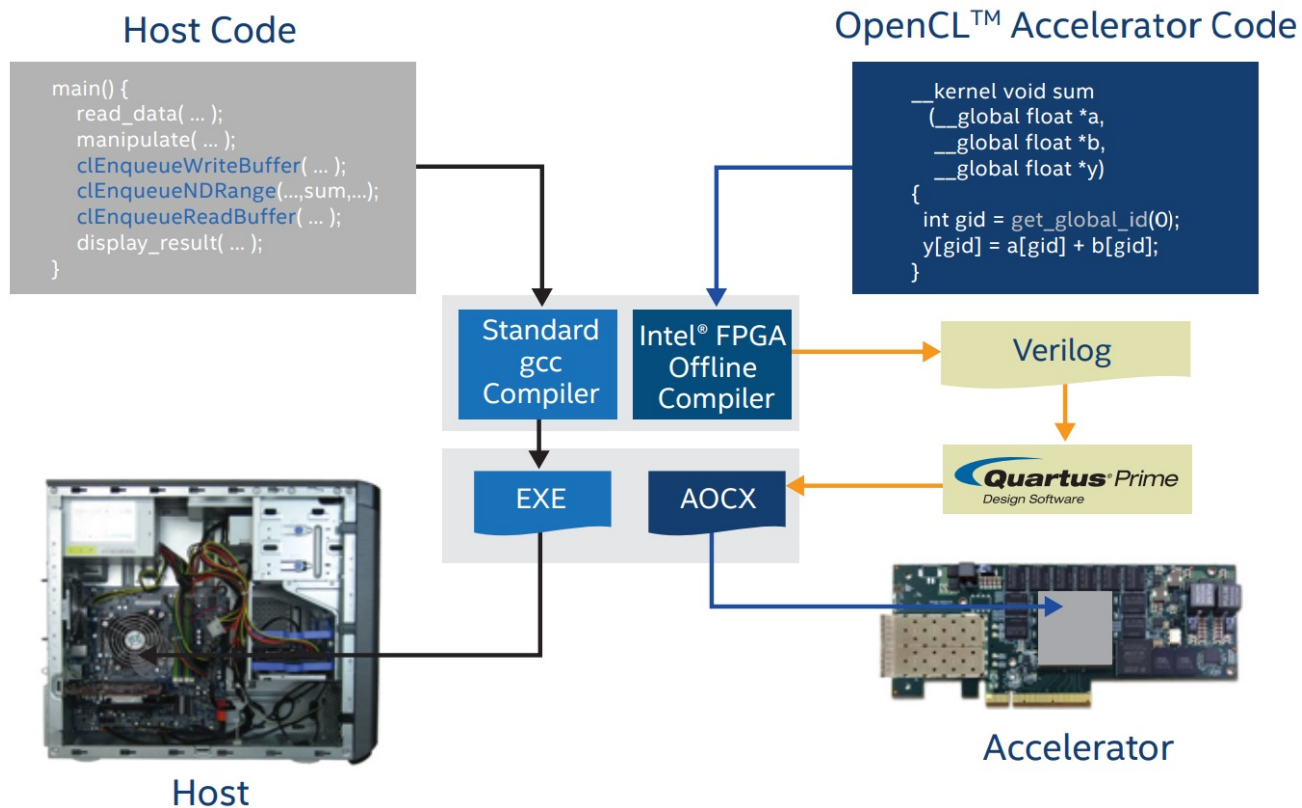
Key challenge: identify the best CPU-FPGA collaboration strategy



Intel Xeon + FPGA Integrated Platform (MCP)

Intel OpenCL SDK for FPGA

- Intel OpenCL SDK for FPGA is used to compile and synthesize host executable and FPGA design



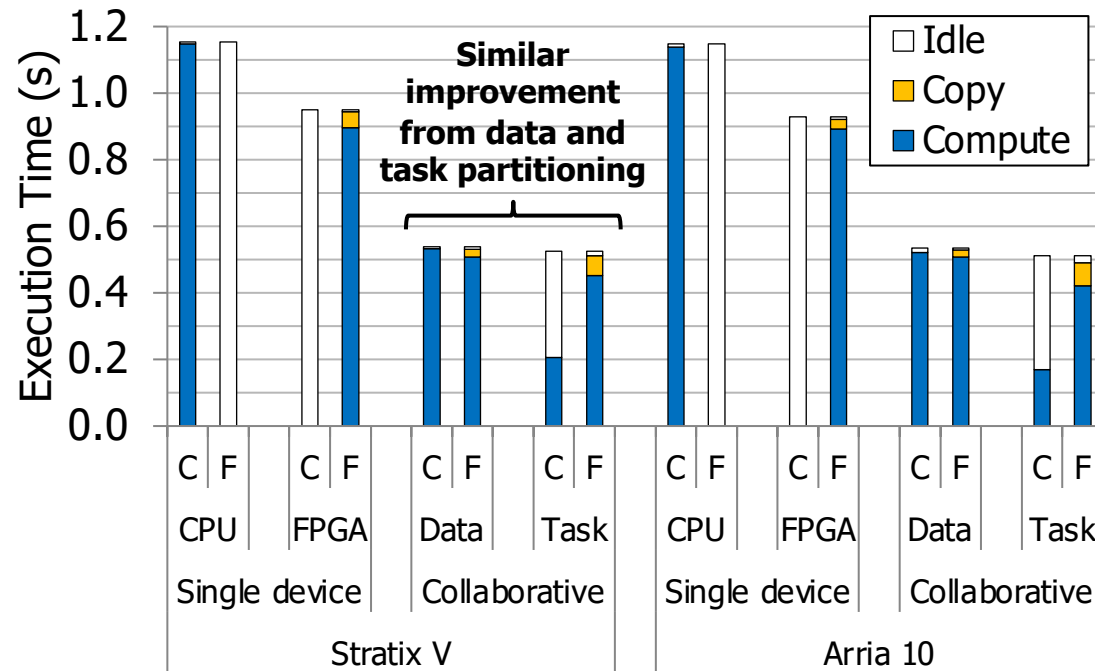
CPU+FPGA Evaluation Platforms



	Platform A	Platform B
FPGA Board	Terasic DE5-Net	Nallatech 510T
FPGA Chip	Intel Stratix V GX	Intel Arria 10 GX
On-Board Memory	4 GB (DDR3)	8 GB (DDR4)
Host CPU	Intel Xeon E3-1240 v3	Intel Xeon E5-2650 v3
Host Memory	8 GB (DDR3)	96 GB (DDR4)
Interface	PCIe gen3.0 x8	PCIe gen3.0 x8

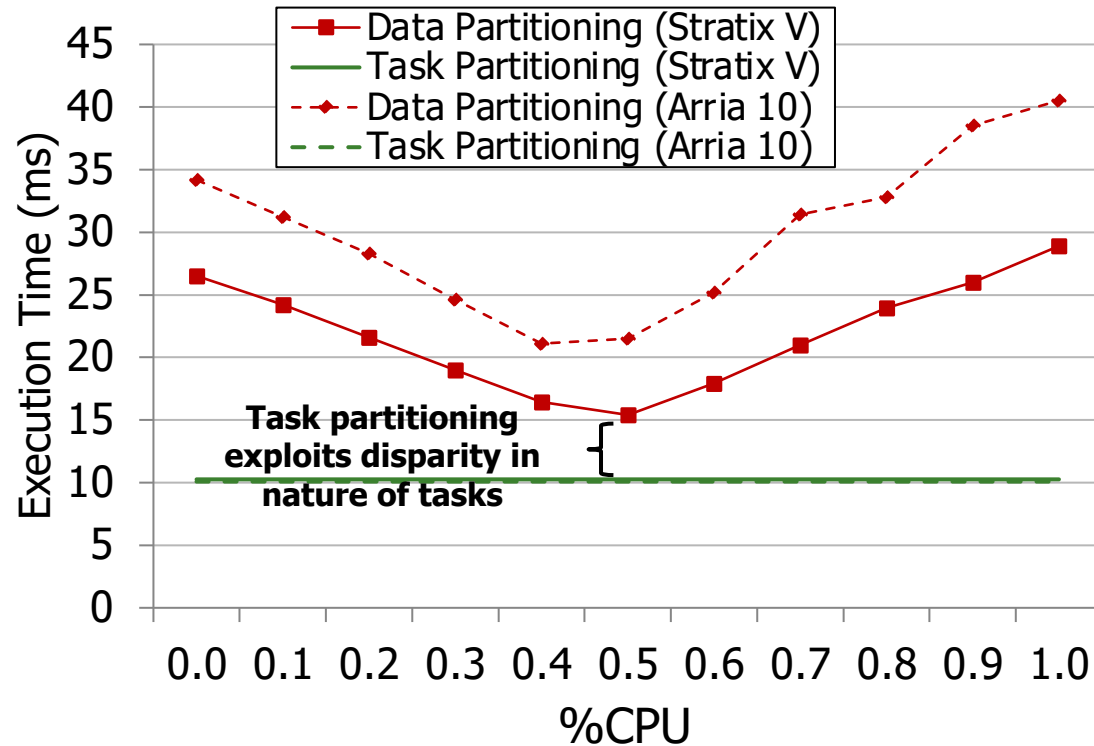
Benefits of Collaboration on FPGA (I)

Case Study: Canny Edge Detection



Benefits of Collaboration on FPGA (II)

Case Study:
Random
Sample
Consensus



Chai on CPU-FPGA Systems (I)

- Sitao Huang, Li-Wen Chang, Izzat El Hajj, Simon Garcia De Gonzalo, Juan Gomez-Luna, Sai Rahul Chalamalasetti, Mohamed El-Hadedy, Dejan Milojicic, Onur Mutlu, Deming Chen, and Wen-mei Hwu,
"Analysis and Modeling of Collaborative Execution Strategies for Heterogeneous CPU-FPGA Architectures"
Proceedings of the 10th ACM/SPEC International Conference on Performance Engineering (ICPE), Mumbai, India, April 2019.
[Slides (pptx)] [pdf]
[Chai CPU-FPGA Benchmark Suite]

Analysis and Modeling of Collaborative Execution Strategies for Heterogeneous CPU-FPGA Architectures

Sitao Huang
ECE, UIUC
shuang91@illinois.edu

Li-Wen Chang*
Microsoft
liwen.chang@microsoft.com

Izzat El Hajj
ECE, UIUC
elhajj2@illinois.edu

Simon Garcia De Gonzalo
CS, UIUC
grcdgnz2@illinois.edu

Juan Gómez-Luna
CS, ETH Zurich
juang@ethz.ch

Sai Rahul Chalamalasetti
Hewlett Packard Labs
sairahul.chalamalasetti@hpe.com

Mohamed El-Hadedy
ECE, Cal Poly Pomona
mealy@cpp.edu

Dejan Milojicic
Hewlett Packard Labs
dejan.milojicic@hpe.com

Onur Mutlu
CS, ETH Zurich
omutlu@ethz.ch

Deming Chen
ECE, UIUC
dchen@illinois.edu

Wen-mei Hwu
ECE, UIUC
w-hwu@illinois.edu

Chai on CPU-FPGA Systems (II)

- Jiantong Jiang, Zeke Wang, Xue Liu, Juan Gómez-Luna, Nan Guan, Qingxu Deng, Wei Zhang, and Onur Mutlu,
"Boyi: A Systematic Framework for Automatically Deciding the Right Execution Model of OpenCL Applications on FPGAs"
Proceedings of the 28th International Symposium on Field-Programmable Gate Arrays (FPGA), Seaside, CA, USA, February 2020.
[[Slides \(pptx\)](#)] [[pdf](#)]

Boyi: A Systematic Framework for Automatically Deciding the Right Execution Model of OpenCL Applications on FPGAs

Jiantong Jiang^{1★}

Nan Guan³

Zeke Wang^{2★}

Qingxu Deng¹

Xue Liu^{1*}

Wei Zhang⁴

Juan Gómez-Luna²

Onur Mutlu²

¹ Department of Computer Science and Engineering, Northeastern University, China

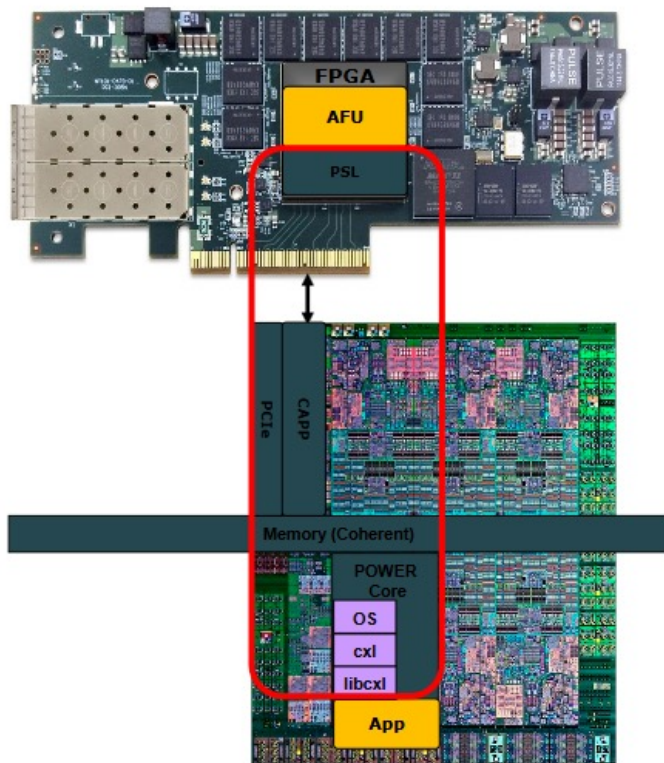
² ETH Zürich, Switzerland

³ Department of Computing, Hong Kong Polytechnic University, Hong Kong

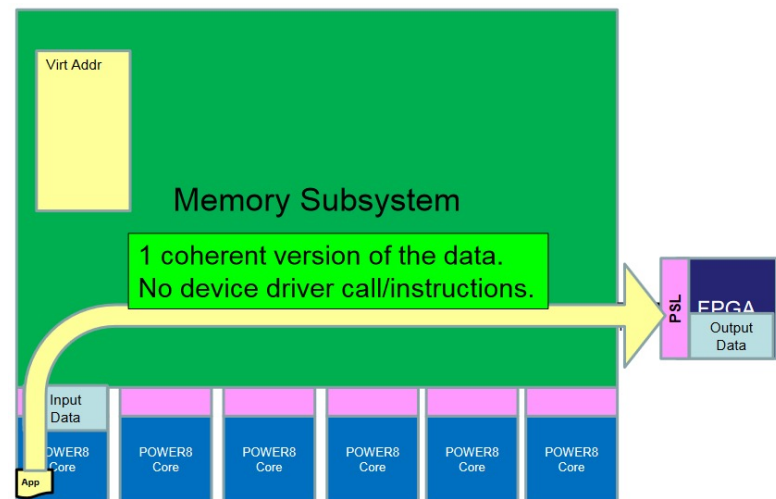
⁴ Department of Electronic and Computer Engineering, Hong Kong University of Science and Technology, Hong Kong

CAPI/OpenCAPI Overview

- CAPI/CAPI2 (Coherent Accelerator Processor Interface)
- OpenCAPI

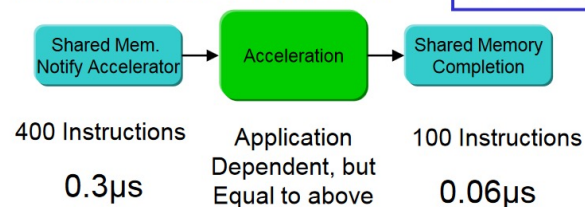


POWER8 - POWER9 Processor



Flow with a CAPI Model:

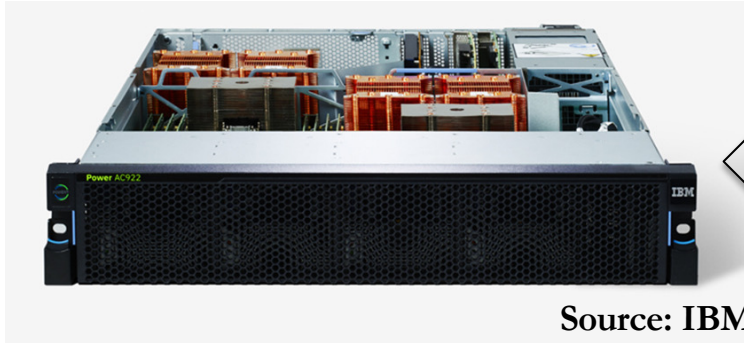
Total 0.36 μ s



Dionysios Diamantopoulos, IBM Research – Zurich, COOL Chips 2018

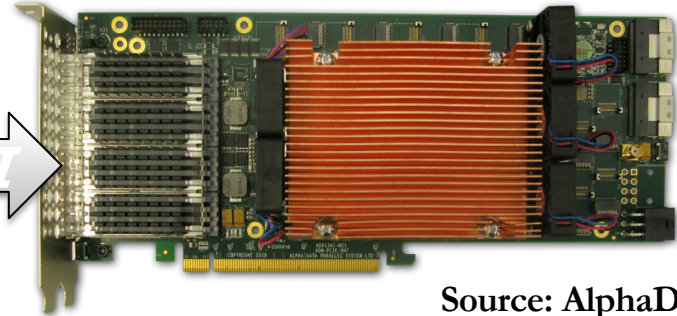


Evaluation Setup for Weather Acceleration



Source: IBM

POWER9 AC922



Source: AlphaData

**HBM-based AD9H7
board**

■ Host System

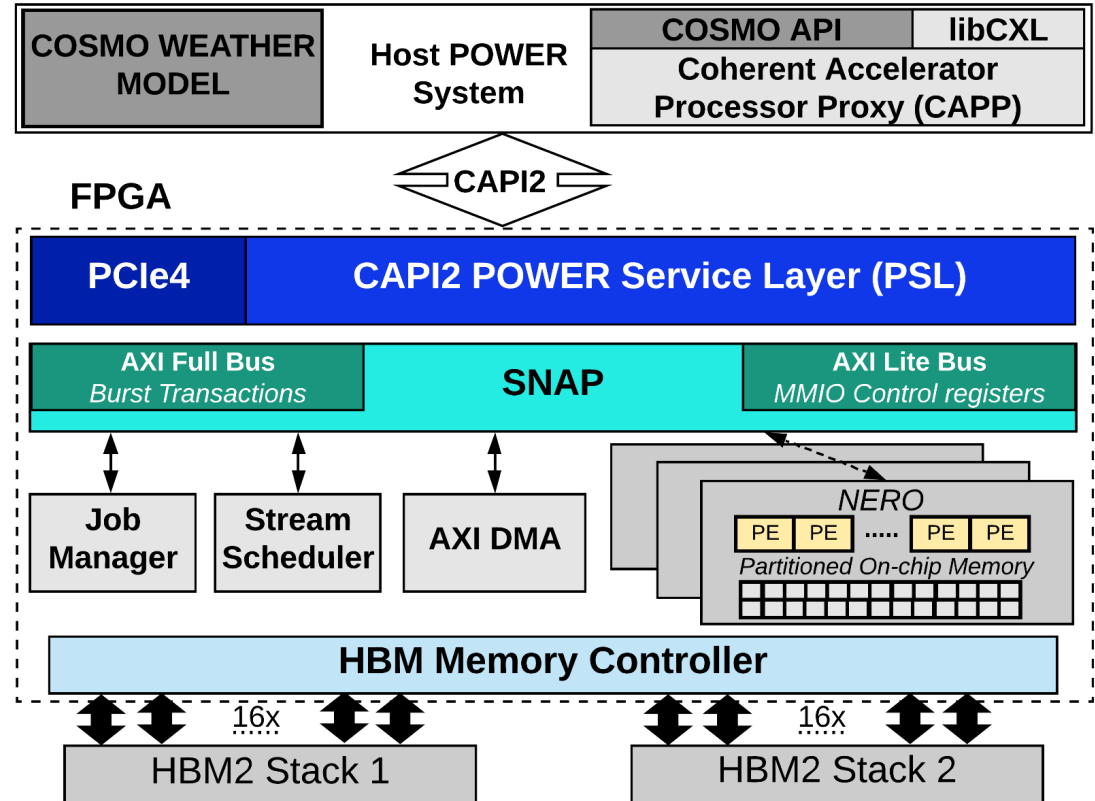
IBM POWER9-16 core (64-threads)

■ FPGA board

Xilinx Virtex® Ultrascale+™ XCVU37P-2

NERO Application Framework

- NERO communicates to Host over **CAPI2** (Coherent Accelerator Processor Interface)
- COSMO API handles offloading jobs to NERO
- SNAP (Storage, Network, and Analytics Programming) allows for seamless integration of the COSMO API



<https://github.com/open-power/snap>

Accelerating Climate Modeling (I)

- Gagandeep Singh, Dionysios Diamantopoulos, Christoph Hagleitner, Juan Gómez-Luna, Sander Stuijk, Onur Mutlu, and Henk Corporaal,
"NERO: A Near High-Bandwidth Memory Stencil Accelerator for Weather Prediction Modeling"

Proceedings of the 30th International Conference on Field-Programmable Logic and Applications (FPL), Gothenburg, Sweden, September 2020.

[[Slides \(pptx\)](#) ([pdf](#))]

[[Lightning Talk Slides \(pptx\)](#) ([pdf](#))]

[[Talk Video](#) (23 minutes)]

Nominated for the Stamatis Vassiliadis Memorial Award.

NERO: A Near High-Bandwidth Memory Stencil Accelerator for Weather Prediction Modeling

Gagandeep Singh^{a,b,c} Dionysios Diamantopoulos^c Christoph Hagleitner^c Juan Gómez-Luna^b

Sander Stuijk^a

Onur Mutlu^b

Henk Corporaal^a

^aEindhoven University of Technology

^bETH Zürich

^cIBM Research Europe, Zurich

Accelerating Climate Modeling (II)

- Gagandeep Singh, Mohammed Alser, Damla Senol Cali, Dionysios Diamantopoulos, Juan Gómez-Luna, Henk Corporaal, and Onur Mutlu, ["FPGA-based Near-Memory Acceleration of Modern Data-Intensive Applications"](#) *IEEE Micro (IEEE MICRO)*, 2021.

FPGA-based Near-Memory Acceleration of Modern Data-Intensive Applications

Gagandeep Singh[◇] Mohammed Alser[◇] Damla Senol Cali[⌘]

Dionysios Diamantopoulos[▽] Juan Gómez-Luna[◇]

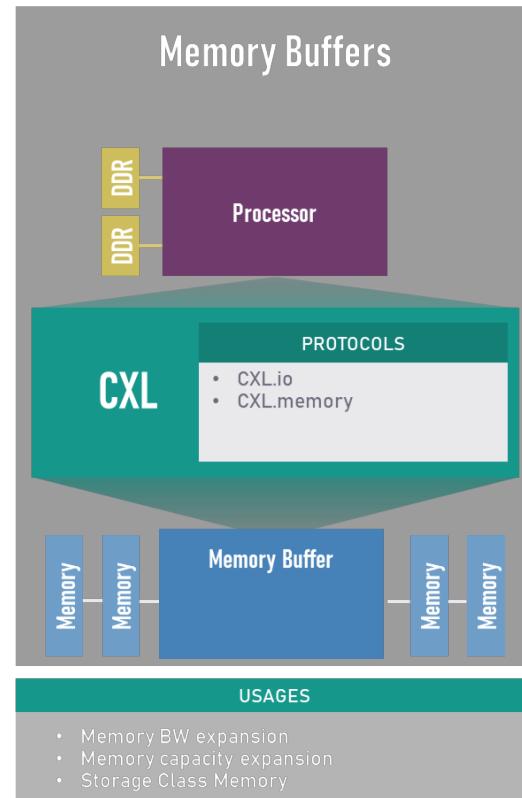
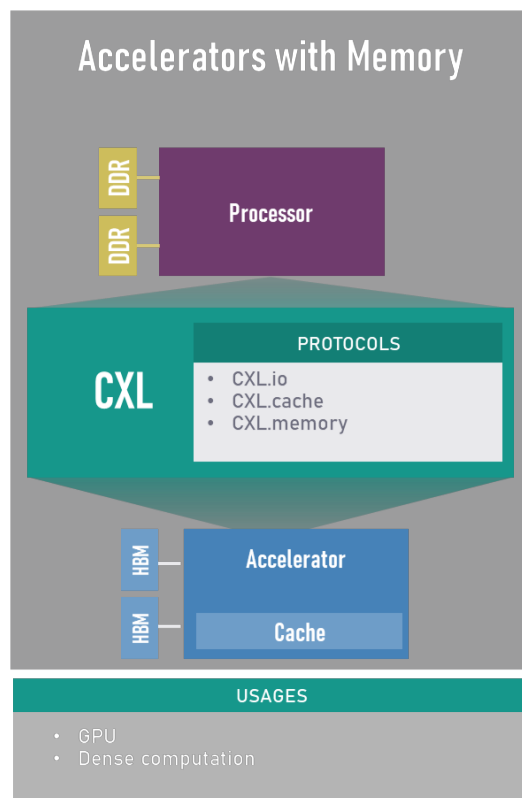
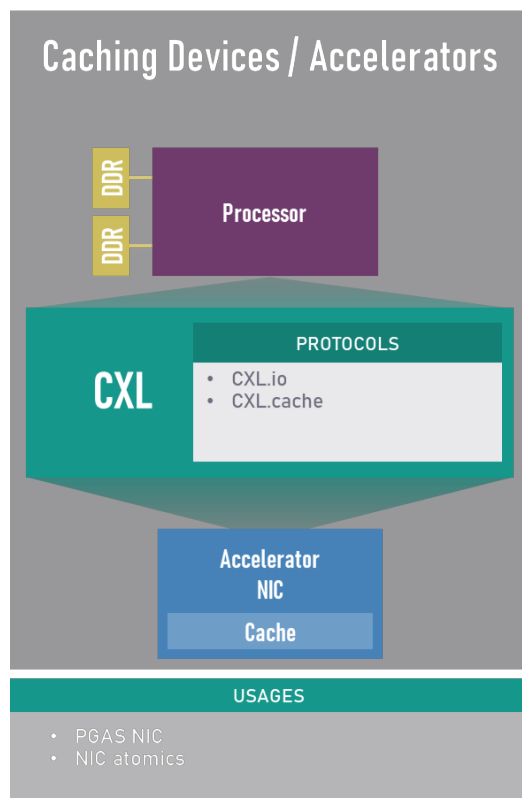
Henk Corporaal^{*} Onur Mutlu^{◇⌘}

[◇]*ETH Zürich* [⌘]*Carnegie Mellon University*

^{*}*Eindhoven University of Technology* [▽]*IBM Research Europe*

Compute Express Link (CXL)

- Compute Express Link (CXL) is an open industry standard interconnect offering **high-bandwidth, low-latency connectivity between host processor and devices** such as accelerators, memory buffers, and smart I/O devices



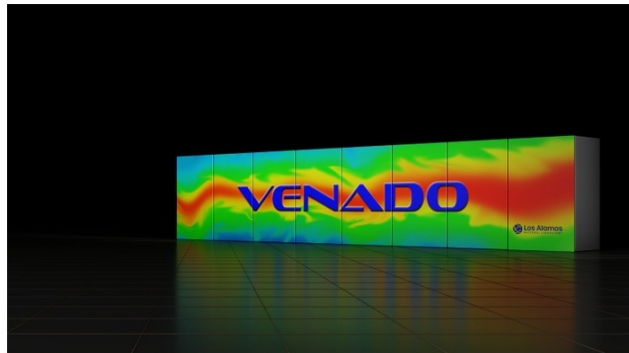
NVIDIA Grace CPU and Grace Hopper



Top Global Systems Makers Accelerate Adoption of NVIDIA Grace and Grace Hopper

Atos, Dell Technologies, GIGABYTE, Hewlett Packard Enterprise, Inspur, Lenovo and Supermicro Join First Wave Planning NVIDIA Grace-Powered HGX Systems for HPC and AI

Monday, May 30, 2022



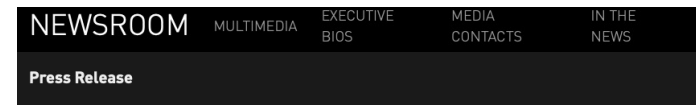
NVIDIA today announced that a range of the world's leading computer makers are adopting the new NVIDIA Grace™ superchips to create the next generation of servers turbocharging AI and HPC workloads for the exascale era.

Atos, Dell Technologies, GIGABYTE, HPE, Inspur, Lenovo and Supermicro are planning to deploy servers built with the [NVIDIA Grace CPU Superchip and NVIDIA Grace Hopper™ Superchip](#).

All these new systems benefit from the just-announced [Grace and Grace Hopper designs in the NVIDIA HGXTM platform](#), which provide manufacturers the blueprints needed to build systems that offer the highest performance and twice the memory bandwidth and energy efficiency of today's leading data center CPU.

"As supercomputing enters the era of exascale AI, NVIDIA is teaming up with our OEM partners to enable researchers to tackle massive challenges previously out of reach," said Ian Buck, vice president of Hyperscale and HPC at NVIDIA. "Across climate science, energy research, space exploration, digital biology, quantum computing and more, the NVIDIA Grace CPU Superchip and Grace Hopper Superchip form the foundation of the world's most advanced platform for HPC and AI."

<https://nvidianews.nvidia.com/news/top-global-systems-makers-accelerate-adoption-of-nvidia-grace-and-grace-hopper>



Taiwan's Tech Titans Adopt World's First NVIDIA Grace CPU-Powered System Designs

New Class of Data Center Systems for Digital Twins, AI, High Performance Computing, Cloud Graphics and Gaming to Come From ASUS, Foxconn Industrial Internet, GIGABYTE, QCT, Supermicro, Wiyynn

Monday, May 23, 2022



COMPUTEX -- NVIDIA today announced that Taiwan's leading computer makers are set to release the first wave of systems powered by the [NVIDIA Grace™ CPU Superchip and Grace Hopper Superchip](#) for a wide range of workloads spanning digital twins, AI, high performance computing, cloud graphics and gaming.

Dozens of server models from [ASUS](#), [Foxconn Industrial Internet](#), [GIGABYTE](#), [QCT](#), [Supermicro](#) and [Wiyynn](#) are expected starting in the first half of 2023. The Grace-powered systems will join x86 and other Arm-based servers to offer customers a broad range of choice for achieving high performance and efficiency in their data centers.

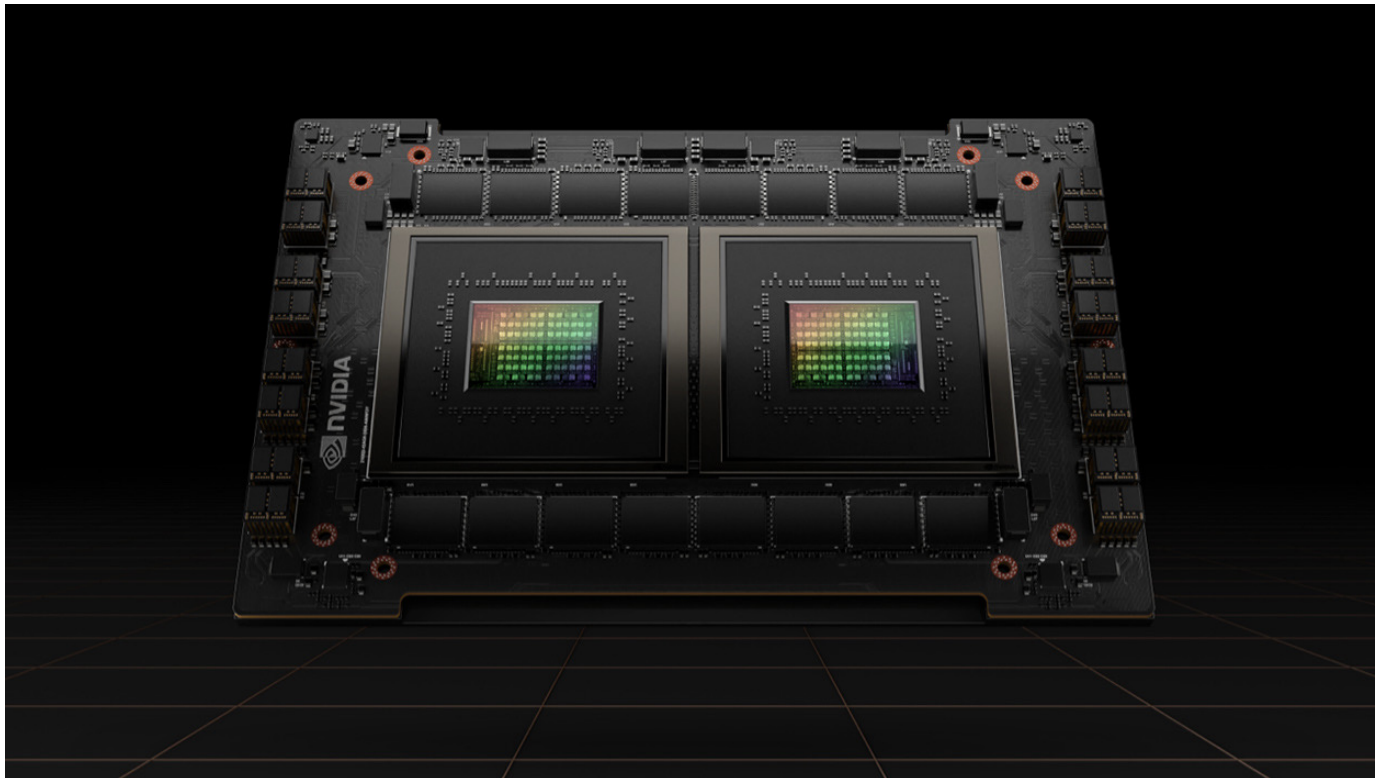
"A new type of data center is emerging — AI factories that process and refine mountains of data to produce intelligence — and NVIDIA is working closely with our Taiwan partners to build the systems that enable this transformation," said Ian Buck, vice president of Hyperscale and HPC at NVIDIA. "These new systems from our partners, powered by our Grace Superchips, will bring the power of accelerated computing to new markets and industries globally."

The coming servers are based on four new system designs featuring the Grace CPU Superchip and Grace Hopper Superchip, which NVIDIA announced at its two most recent GTC conferences. The 2U form factor designs provide the blueprints and server baseboards for original design manufacturers and original equipment manufacturers to quickly bring to market systems for the NVIDIA CGX™ cloud gaming, [NVIDIA OVX™](#) digital twin and the [NVIDIA HGX™](#) AI and HPC platforms.

<https://nvidianews.nvidia.com/news/taiwans-tech-titans-adopt-worlds-first-nvidia-grace-cpu-powered-system-designs>

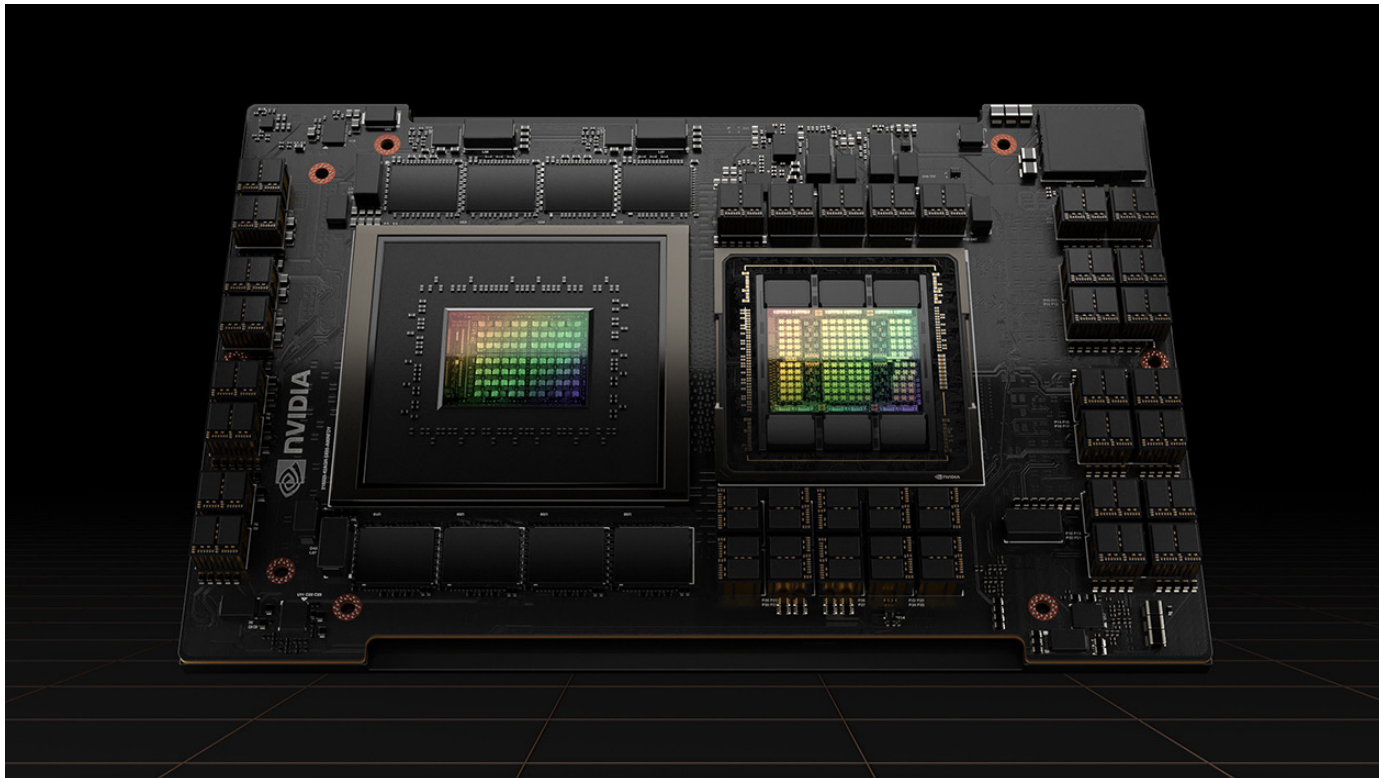
NVIDIA Grace CPU Superchip

- 144 ARM v9 CPU cores
- LPDDR5x memory with ECC, 1 TB/s total bandwidth
- 900 GB/s coherent interface (7x faster than PCIe Gen 5)



NVIDIA Grace Hopper Superchip

- CPU + GPU
 - Grace CPU + Hopper GPU
- 900 GB/s coherent interface (7x faster than PCIe Gen 5)

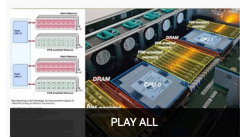


Collaborative Computing: Key Takeaways

- Possibility of having **several devices collaborating** on the same workload
- And having **the most appropriate cores** for each workload, exploiting heterogeneity
- Easier programming with Unified Memory or Shared Virtual Memory
- CPU-GPU memory coherence and system-wide atomic operations since NVIDIA Pascal and HSA
 - Fine-grain collaboration

Processing-in-Memory Course (Spring 2022)

- Short weekly lectures
- Hands-on projects



Livestream - P&S Exploring the Processing-in-Memory Paradigm for Future Computing Systems (Spring 2022)

14 videos • 580 views • Updated 6 days ago



- 1 Processing-in-Memory Course: Lecture 1: Exploring the PIM Paradigm for Future Systems - Spring 2022
Onur Mutlu Lectures
1:35:48
- 2 Processing-in-Memory Course: Lecture 2: Real-world PIM: UPMEM PIM Architecture - Spring 2022
Onur Mutlu Lectures
31:53
- 3 Processing-in-Memory Course: Lecture 3: Real-world PIM: Microbenchmarking of UPMEM PIM - Spring 2022
Onur Mutlu Lectures
56:51
- 4 Processing-in-Memory Course: Lecture 4: Real-world PIM: Samsung HBM-PIM Architecture - Spring 2022
Onur Mutlu Lectures
1:06:06
- 5 Processing-in-Memory Course: Lecture 5: How to Evaluate Data Movement Bottlenecks - Spring 2022
Onur Mutlu Lectures
1:00:04
- 6 Processing-in-Memory Course: Lecture 6: Real-world PIM: SK Hynix AiM - Spring 2022
Onur Mutlu Lectures
45:00
- 7 Processing-in-Memory Course: Lecture 7: Programming PIM Architectures - Spring 2022
Onur Mutlu Lectures
46:35
- 8 Processing-in-Memory Course: Lecture 8: Benchmarking and Workload Suitability on PIM - Spring 2022
Onur Mutlu Lectures
34:36
- 9 Processing-in-Memory Course: Lecture 9: Real-world PIM: Samsung AxDIMM - Spring 2022
Onur Mutlu Lectures
32:37

https://youtube.com/playlist?list=PL5Q2soXY2Zi-0NK1C5vi2Zx9nmE_3-cKN



SAFARI Project & Seminars Courses
(Spring 2022)

Trace: • heterogeneous_systems • processing_in_memory

Home

Courses

- SoftMC
- Ramulator
- Accelerating Genomics
- Mobile Genomics
- **Processing-in-Memory**
- Heterogeneous Systems
- Modern SSDs

Exploring the Processing-in-Memory Paradigm for Future Computing Systems

Course Description

Data movement between the memory units and the compute units of current computing systems is a major performance and energy bottleneck. From large-scale servers to mobile devices, data movement costs dominate computation costs in terms of both performance and energy consumption. For example, data movement between the main memory and the processing cores accounts for 62% of the total system energy in consumer applications. As a result, the data movement bottleneck is a huge burden that greatly limits the energy efficiency and performance of modern computing systems. This phenomenon is an undesired effect of the dichotomy between memory and the processor, which leads to the data movement bottleneck.

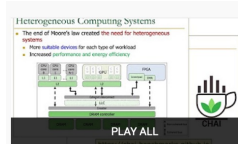
Many modern and important workloads such as machine learning, computational biology, graph processing, databases, video analytics, and real-time data analytics suffer greatly from the data movement bottleneck. These workloads are exemplified by irregular memory accesses, relatively low data reuse, low cache line utilization, low arithmetic intensity (i.e., ratio of operations per accessed byte), and large datasets that greatly exceed the main memory size. The computation in these workloads cannot usually compensate for the data movement costs. In order to alleviate this data movement bottleneck, we need a paradigm shift from the traditional processor-centric design, where all computation takes place in the compute units, to a more data-centric design where processing elements are placed closer to or inside where the data resides. This paradigm of computing is known as Processing-in-Memory (PIM).

This is your perfect P&S if you want to become familiar with the main PIM technologies, which represent "the next big thing" in Computer Architecture. You will work hands-on with the first real-world PIM architecture, will explore different PIM architecture designs for important workloads, and will develop tools to enable research of future PIM systems. Projects in this course span software and hardware as well as the software/hardware interface. You can potentially work on developing and optimizing new workloads for the first real-world PIM hardware or explore new PIM designs in simulators, or do something else that can forward our understanding of the PIM paradigm.

https://safari.ethz.ch/projects_and_seminars/spring2022/doku.php?id=processing_in_memory

Heterogeneous Systems Course (Spring 2022)

- Short weekly lectures
- Hands-on projects



Livestream - P&S Hands-on Acceleration on Heterogeneous Computing Systems (Spring 2022)

13 videos • 889 views • Updated 6 days ago



Onur Mutlu Lectures

SUBSCRIBED



10 videos in playlist:

- HetSys Course: Lecture 1: Hands-on Acceleration on Heterogeneous Computing Systems (Spring 2022) - Onur Mutlu Lectures - 41:54
- HetSys Course: Lecture 2: SIMD Processing and GPUs (Spring 2022) - Onur Mutlu Lectures - 1:22:48
- HetSys Course: Lecture 3: GPU Software Hierarchy (Spring 2022) - Onur Mutlu Lectures - 56:24
- HetSys Course: Lecture 4: GPU Memory Hierarchy (Spring 2022) - Onur Mutlu Lectures - 54:27
- HetSys Course: Lecture 5: GPU Performance Considerations (Spring 2022) - Onur Mutlu Lectures - 1:23:29
- HetSys Course: Lecture 6: Parallel Patterns: Reduction (Spring 2022) - Onur Mutlu Lectures - 33:39
- HetSys Course: Lecture 7: Parallel Patterns: Histogram (Spring 2022) - Onur Mutlu Lectures - 1:29:40
- HetSys Course: Lecture 8: Parallel Patterns: Convolution (Spring 2022) - Onur Mutlu Lectures - 1:03:15
- HetSys Course: Lecture 9: Parallel Patterns: Prefix Sum (Scan) (Spring 2022) - Onur Mutlu Lectures - 1:19:46
- HetSys Course: Lecture 10: Parallel Patterns: Sparse Matrices (Spring 2022) - Onur Mutlu Lectures - 1:19:46

SAFARI Project & Seminars Courses (Spring 2022)

Trace: • processing_in_memory • heterogeneous_systems

Home

Courses

- SoftMC
- Ramulator
- Accelerating Genomics
- Mobile Genomics
- Processing-in-Memory

heterogeneous_systems

Hands-on Acceleration on Heterogeneous Computing Systems

Course Description

The increasing difficulty of scaling the performance and efficiency of CPUs every year has created the need for turning computers into heterogeneous systems, i.e., systems composed of multiple types of processors that can suit better different types of workloads or parts of them. More than a decade ago, Graphics Processing Units (GPUs) became general-purpose parallel processors, in order to make their outstanding processing capabilities available to many workloads beyond graphics. GPUs have been critical key to the recent rise of Machine Learning and Artificial Intelligence, which took unrealistic training times before the use of GPUs. Field-Programmable Gate Arrays (FPGAs) are another example computing device that can deliver impressive benefits in terms of performance and energy efficiency. More specific examples are (1) a plethora of specialized accelerators (e.g., Tensor Processing Units for neural networks), and (2) near-data processing architectures (i.e., placing compute capabilities near or inside memory/storage).

Despite the great advances in the adoption of heterogeneous systems in recent years, there are still many challenges to tackle, for example:

- Heterogeneous implementations (using GPUs, FPGAs, TPUs) of modern applications from important fields such as bioinformatics, machine learning, graph processing, medical imaging, personalized medicine, robotics, virtual reality, etc.
- Scheduling techniques for heterogeneous systems with different general-purpose processors and accelerators, e.g., kernel offloading, memory scheduling, etc.
- Workload characterization and programming tools that enable easier and more efficient use of heterogeneous systems.

If you are enthusiastic about working **hands-on** with different software, hardware, and architecture projects for heterogeneous systems, this is your P&S. You will have the opportunity to program heterogeneous systems with different types of devices (CPUs, GPUs, FPGAs, TPUs), propose algorithmic changes to important applications to better leverage the compute power of heterogeneous systems, understand different workloads and identify the most suitable device for their execution, design optimized scheduling techniques, etc. In general, the goal will be to reach the highest performance reported for a given important application.

Table of Contents

- Hands-on Acceleration on Heterogeneous Computing Systems
- Course Description
- Mentors
- Lecture Video Playlists on YouTube
- Spring 2022 Meetings/Schedule
- Learning Materials
- Assignments

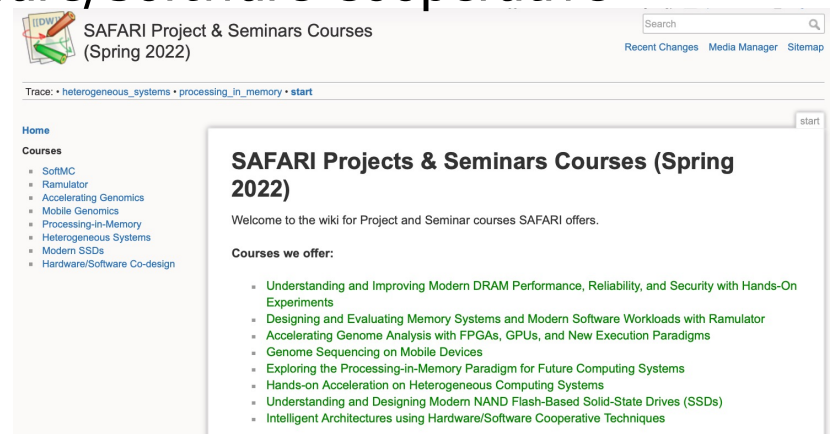
https://youtube.com/playlist?list=PL5Q2soXY2Zi9XrqXR38IM_FTjmY6h7Gzm

https://safari.ethz.ch/projects_and_seminars/spring2022/doku.php?id=heterogeneous_systems

More P&S Courses: SSDs, Memory, Bioinformatics...

- Understanding and Improving Modern DRAM Performance, Reliability, and Security with Hands-On Experiments
- Designing and Evaluating Memory Systems and Modern Software Workloads with Ramulator
- Accelerating Genome Analysis with FPGAs, GPUs, and New Execution Paradigms
- Genome Sequencing on Mobile Devices
- Understanding and Designing Modern NAND Flash-Based Solid-State Drives (SSDs)
- Intelligent Architectures using Hardware/Software Cooperative Techniques

https://safari.ethz.ch/projects_and_seminars/spring2022/doku.php?id=start



More Resources: Onur Mutlu Lectures

- All P&S courses
- Digital Design and CompArch course
- Advanced CompArch course
- Seminar in CompArch

The screenshot displays the YouTube channel page for "Onur Mutlu Lectures", which has 25.6K subscribers. The channel is subscribed to. The "PLAYLISTS" tab is selected, showing a list of created playlists. The first section, "Created playlists", includes:

- SAFARI Live Seminars 2022**: Updated 4 days ago. View Full Playlist.
- Livestream - P&S Modern SSDs** (Spring 2022): Updated 6 days ago. View Full Playlist.
- Livestream - P&S Intelligent Architectures via...**: Updated 6 days ago. View Full Playlist.
- Livestream - P&S Hands-on Acceleration on Heterogeneo...**: Updated 6 days ago. View Full Playlist.
- Livestream - P&S Exploring the Processing-in-Memory...**: Updated 6 days ago. View Full Playlist.
- Livestream - P&S Accelerating Genome Analysis with FPGAs...**: Updated 7 days ago. View Full Playlist.

The second section, "First Course in Computer Architecture & Digital Design 2021-2013", includes:

- Livestream - Digital Design and Computer Architecture - ETH...**: Onur Mutlu Lectures. View Full Playlist.
- Digital Design & Computer Architecture - ETH Zürich...**: Onur Mutlu Lectures. View Full Playlist.
- Design of Digital Circuits - ETH Zürich - Spring 2019**: Onur Mutlu Lectures. View Full Playlist.
- Design of Digital Circuits - ETH Zürich - Spring 2018**: Onur Mutlu Lectures. View Full Playlist.
- Digital Circuits and Computer Architecture - ETH Zürich - ...**: Onur Mutlu Lectures. View Full Playlist.
- Spring 2015 - Computer Architecture Lectures - ...**: Carnegie Mellon Computer Architecture. View Full Playlist.

The third section, "Advanced Computer Architecture Courses 2021-2012", includes:

- Processing in GEMM Engine**: Includes standard GEMM modules, number of CPU processors, and large amounts of context.
- Look back to the past**: Includes the design and implementation of the GEMM engine, the design of the GEMM engine, and the design of the GEMM engine.
- Look back to the past**: Includes the design and implementation of the GEMM engine, the design of the GEMM engine, and the design of the GEMM engine.

P&S Heterogeneous Systems

Collaborative Computing

Dr. Juan Gómez Luna

Prof. Onur Mutlu

ETH Zürich

Fall 2022

16 January 2023